

# Explicit cost bounds of stochastic Galerkin approximations for parameterized PDEs with random coefficients<sup>☆</sup>

N.C. Dexter<sup>a</sup>, C.G. Webster<sup>b</sup>, G. Zhang<sup>b</sup>

<sup>a</sup>*Department of Mathematics, University of Tennessee, Knoxville, TN 37996.*

<sup>b</sup>*Department of Computational and Applied Mathematics, Oak Ridge National Laboratory, Oak Ridge, TN 37831.*

---

## Abstract

This work analyzes the overall computational complexity of the stochastic Galerkin finite element method (SGFEM) for approximating the solution of parameterized elliptic partial differential equations with both affine and non-affine random coefficients. To compute the fully discrete solution, such approaches employ a Galerkin projection in both the deterministic and stochastic domains, produced here by a combination of finite elements and a global orthogonal basis, defined on an isotopic total degree index set, respectively. To account for the sparsity of the resulting system, we present a rigorous cost analysis that considers the total number of coupled finite element systems that must be simultaneously solved in the SGFEM. However, to maintain sparsity as the coefficient becomes increasingly nonlinear in the parameterization, it is necessary to also approximate the coefficient by an additional orthogonal expansion. In this case we prove a rigorous complexity estimate for the number of floating point operations (FLOPs) required per matrix-vector multiplication of the coupled system. Based on such complexity estimates we also develop explicit cost bounds in terms of FLOPs to solve the stochastic Galerkin (SG) systems to a prescribed tolerance, which are used to compare with the minimal complexity estimates of a stochastic collocation finite element method (SCFEM), shown in our previous work [16]. Finally, computational evidence complements the theoretical estimates and supports our conclusion that, in the case that the coefficient is affine, the coupled SG system can be solved more efficiently than the decoupled SC systems. However, as the coefficient becomes more nonlinear, it becomes prohibitively expensive to obtain an approximation with the SGFEM.

**Keywords:** stochastic Galerkin, stochastic collocation, sparse polynomial approximation, complexity analysis, explicit cost bounds, finite elements

---

## 1. Introduction

Nowadays, stochastic polynomial methods are widely used alternatives to Monte Carlo methods (see, e.g., [15]) for predicting the solution to physical and engineering problems described by parameterized partial differential equations (PDEs) with a finite number of random variables. In the last decade, two classes of such methods have been proposed that often feature much faster convergence rates: *intrusive* stochastic Galerkin (SG) methods and *non-intrusive* stochastic collocation (SC) methods. Both approaches typically employ a Galerkin projection in the physical domain, produced here by finite elements, and the resulting fully discrete approximations only differ in their choice of multivariate polynomials for the discretization in the stochastic domain. For details about the relations between these methods see [17, 18, 19, 21, 24], and for computational comparisons between the SG and SC methods see, e.g., [3, 13, 19].

---

<sup>☆</sup>This material is based upon work supported in part by the U.S. Air Force of Scientific Research under grant number 1854-V521-12; by the U.S. Department of Energy, Office of Science, Office of Advanced Scientific Computing Research, Applied Mathematics program under contract numbers ERKJ259, and ERKJE45 and by the Laboratory Directed Research and Development program at the Oak Ridge National Laboratory, which is operated by UT-Battelle, LLC, for the U.S. Department of Energy under Contract DE-AC05-00OR22725.

The focus of this paper is to provide explicit cost bounds for applying the stochastic Galerkin finite element method (SGFEM) to the solution of an elliptic PDE, with stochastic diffusion coefficient parameterized by finitely many random variables. In particular, we focus on the cost of constructing isotropic total degree SG approximations when the coefficient has both affine and non-affine dependence on the parameters. Under very basic assumptions on the coefficient, the solution to this problem has been shown to have analytic regularity in the random variables (see [33]). As a result, SG approximations that employ a global orthogonal basis have been shown to be optimal projections in the  $L^2$  sense, converging sub-exponentially with respect to the cardinality of the polynomial subspace [32]. However, the computational cost of solving the coupled SG system does not grow linearly in the cardinality of the given subspace. Therefore, the convergence estimates do not indicate the total complexity of obtaining the approximation for a prescribed tolerance.

When the diffusion coefficient can be written as a sum of separable functions of the physical and random parameters, the coupled SG system can be written as a sum of Kronecker products of SG matrices and finite element stiffness matrices. For every SG matrix, each nonzero element leads to a nonzero block of the coupled SG system, where the size of the block equals the size of the finite element stiffness matrix. To solve the SG system, one must simultaneously solve all the coupled finite element problems. In the case that the coefficient is affine in the parameters, the number of nonzeros in each SG matrix is of order  $\mathcal{O}(M_p)$  [14], where  $M_p$  is the cardinality of the isotropic total degree polynomial subspace of order  $p \in \mathbb{N}$ . Thus, a matrix-vector product involving the coupled SG system requires  $\mathcal{O}(J_h M_p)$  floating point operations (FLOPs), where  $J_h$  is the number of physical degrees of freedom. Therefore, the work of solving the coupled SG system when employing an iterative method, e.g., conjugate gradient (CG), is of the order  $\mathcal{O}(J_h M_p N_{\text{iter}}^{\text{SG}})$  where  $N_{\text{iter}}^{\text{SG}}$  is the number of iterations required to achieve a prescribed accuracy of the fully discrete approximation [3, 14, 34].

On the other hand, when the diffusion coefficient is a general non-affine function of the random parameters, the cost of obtaining an approximation with the SGFEM is not as obvious as before. In this setting we consider two cases, namely, the coefficient is: (1) a polynomial with respect to the random variables, and; (2) a transcendental function with respect to the random variables. In the first case, as we increase the order of the polynomial, the block-sparsity of the SG system decreases, resulting in a SG system that *incrementally* becomes block-dense [12, 14, 21, 23, 34, 35]. In the second case, a separable representation can be guaranteed with the use of an orthogonal expansion [36, 37], such that, substituting the expansion into the discretized PDE recovers the Kronecker product structure. However, when the expansion is not truncated, the SG system is known to be *entirely* block-dense [14, 23]. Without a priori knowledge on the exact sparsity of the SG matrices in this case, it was estimated that the complexity of matrix-vector multiplications of the SG Kronecker product system is between  $\mathcal{O}(J_h M_p^2)$  and  $\mathcal{O}(J_h M_p^3)$  [34]. As such, it is impossible to make a conclusive statement about the computational cost, and, more importantly, does not account for the two cases above, i.e., when the coefficient is possibly a truncated polynomial of fixed total degree  $r \in \mathbb{N}$  such that  $1 \leq r < \infty$ . In these cases, the work of solving the coupled SG system with an iterative method is given by  $\mathcal{O}(J_h \mathcal{M}(p, r) N_{\text{iter}}^{\text{SG}})$ , where  $\mathcal{M}(p, r)$  is the total number of  $\mathcal{O}(J_h)$  finite element problems that must be simultaneously solved.

The key challenge of estimating the cost of solving the SG system when the coefficient is a (truncated) polynomial of finite order is to provide bounds on the block-sparsity of the matrix, i.e., nonzeros of the SG system. To achieve this, we provide a rigorous counting argument, which can be seen as a generalization of results from [14], for the exact sparsity of the SG matrices for an arbitrary order orthogonal expansion of a non-affine coefficient. As a result, we are able to provide bounds for  $\mathcal{M}(p, r)$  of the order  $\mathcal{O}(M_p M_r \min\{2^r, M_{\lceil r/2 \rceil}\})$ , where  $M_r$  is the cardinality of the total degree polynomial subspace used in an orthogonal expansion of order  $r$  of the coefficient. This result provides sharper estimates than the bounds in the case of the full orthogonal expansion from [35] since it depends on the truncation order  $r$ , and allows us to estimate the total complexity of solving the coupled system for general non-affine coefficients. Since the counting argument for the sparsity of the SG system relies only on the SG discretization of an elliptic operator in terms of orthogonal polynomials, we note that this argument can be reused to estimate the complexity of solving similarly defined PDEs with this method.

In addition, we also develop explicit cost bounds in terms of FLOPs to solve the SG system. Our

approach relies on  $\varepsilon$ -complexity analysis, wherein we balance the errors arising from the approximation with the SGFEM and the iterative solver, e.g., CG, so as to ensure the solution to the fully discrete approximation achieves a given tolerance of  $\varepsilon > 0$ . With this result, we are able to provide a direct comparison with  $\varepsilon$ -complexity estimates for the stochastic collocation finite element method (SCFEM) in our previous work [16]. Finally, we present numerical results in agreement with the theoretical work estimates for both the SGFEM and SCFEM all cases described above.

An outline of the paper is as follows. In §2, we provide a discussion on the model problem, and requirements on the diffusion coefficient. In §3, we define the parameterized finite element and SG approximations, derive the SG system, and provide examples of the resulting linear systems that arise from the SG discretization with various coefficients. We then define the cost of solving the SG system and discuss preconditioning strategies. In §4, we derive the exact number of coupled finite element problems in the SG system and bounds on the sparsity in the non-affine case, and present explicit cost bounds of the SGFEM. We also discuss the conditioning of the system in the non-affine case in order to provide a comparison with similar results from [29]. In §5, we briefly describe the SCFEM, and provide theoretical comparison with results from [16] in terms of minimum work to reach a given tolerance, both in the affine and non-affine cases. Finally, in §6, we present illustrative numerical examples corroborating our theoretical results.

## 2. Problem setting

We consider the simultaneous solution of the parameterized linear elliptic PDE:

$$\begin{cases} -\nabla \cdot (a(x, \mathbf{y}) \nabla u(x, \mathbf{y})) = f(x) & \forall x \in D, \mathbf{y} \in \Gamma \\ u(x, \mathbf{y}) = 0 & \forall x \in \partial D, \mathbf{y} \in \Gamma \end{cases} \quad (1)$$

where  $f \in L^2(D)$  is a fixed function of  $x$ ,  $D \subset \mathbb{R}^d$ ,  $d = 1, 2, 3$ , is a bounded Lipschitz domain, and  $\mathbf{y}(\omega) = (y_1(\omega), \dots, y_N(\omega)) : \Omega \rightarrow \Gamma = \prod_{i=1}^N \Gamma_i \subseteq \mathbb{R}^N$  is a random vector with  $\omega \in \Omega$  and  $\Omega$  the set of outcomes. In this setting we assume the components of  $\mathbf{y}$  have a joint probability density function  $\varrho : \Gamma \rightarrow \mathbb{R}_+$ , with  $\varrho(\mathbf{y}) = \prod_{i=1}^N \varrho_i(y_i)$  known directly through, e.g., truncations of correlated random fields [22] in  $(\Gamma, \mathcal{B}(\Gamma), \varrho(\mathbf{y})d\mathbf{y})$ , where  $\mathcal{B}(\mathbf{y})$  denotes the Borel  $\sigma$ -algebra on  $\Gamma$  and  $\varrho(\mathbf{y})d\mathbf{y}$  is the probability measure of  $\mathbf{y}$ . We further assume that  $\varrho_i$  is an even weight function for each  $i = 1, \dots, N$ . We require the following assumptions related to the continuity, coercivity, and holomorphic dependence of the coefficient  $a(x, \mathbf{y})$ . Namely:

- (A1) *There exist constants  $0 < a_{\min} \leq a_{\max} < \infty$  such that for all  $x \in \overline{D}$  and  $\mathbf{y} \in \Gamma$ ,  $a_{\min} \leq a(x, \mathbf{y}) \leq a_{\max}$ .*
- (A2) *The complex continuation of  $a(x, \mathbf{y})$ , denoted  $a^* : \mathbb{C}^N \rightarrow L^\infty$ , is a  $L^\infty(D)$ -valued holomorphic function on  $\mathbb{C}^N$ .*

The holomorphic dependence on  $\mathbf{y}$  of the coefficient  $a(x, \mathbf{y})$  holds in many examples, including polynomial, exponential, and trigonometric functions of the variables  $y_1, \dots, y_N$  shown below.

**Example 2.1** (The affine case). *We consider an affine function of the random parameters, e.g.,*

$$a(x, \mathbf{y}) = a_0(x) + \sum_{k=1}^N y_k b_k(x), \quad x \in \overline{D}, \mathbf{y} \in \Gamma, \quad (2)$$

where  $a_0, \{b_k\}_{k=1}^N \subset L^2(D)$  are such that  $a(x, \mathbf{y})$  satisfies (A1). Such examples include general Karhunen-Loève expansions [22] or piecewise constant random fields.

**Example 2.2** (The non-affine, polynomial case). *We consider a non-affine, polynomial function of the random parameters, e.g.,*

$$a(x, \mathbf{y}) = a_0(x) + \sum_{1 \leq |\alpha| \leq \bar{r}} \mathbf{y}^\alpha c_\alpha(x), \quad x \in \overline{D}, \mathbf{y} \in \Gamma, \quad (3)$$

where  $\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_N)$  is a multi-index,  $|\boldsymbol{\alpha}| = \alpha_1 + \dots + \alpha_N$ ,  $\mathbf{y}^{\boldsymbol{\alpha}} = y_1^{\alpha_1} \dots y_N^{\alpha_N}$ ,  $\bar{r} < \infty$  is the polynomial order of  $a(x, \mathbf{y})$ , and  $a_0, \{c_{\boldsymbol{\alpha}}\}_{|\boldsymbol{\alpha}| \leq \bar{r}} \subset L^2(D)$  are such that  $a(x, \mathbf{y})$  satisfies (A1). Examples include fixed-order Taylor or orthogonal expansions of a general random field.

**Example 2.3** (The non-affine, transcendental case). We consider a non-affine, transcendental function of the random parameters, e.g.,

$$a(x, \mathbf{y}) = a_0(x) + g(x, \mathbf{y}), \quad x \in \bar{D}, \mathbf{y} \in \Gamma, \quad (4)$$

where  $a_0, g \subset L^2(D)$ , and  $g(x, \mathbf{y})$  is a general transcendental function of  $x$  and  $\mathbf{y}$ , such that  $a(x, \mathbf{y})$  satisfies (A1). Examples of  $g(x, \mathbf{y})$  include the sine, logarithm, or exponential functions of (2) or (3).

Let  $L^2_{\varrho}(\Gamma)$  be the space of square integrable functions with respect to the measure  $\varrho(\mathbf{y})d\mathbf{y}$  and  $L^{\infty}_{\varrho}(\Gamma)$  be the space of essentially bounded functions, with the norm

$$\|u\|_{L^{\infty}_{\varrho}(\Gamma)} := \operatorname{ess\,sup}_{\mathbf{y} \in \Gamma} |u(\mathbf{y})|,$$

where the essential supremum is taken with respect to the weight  $\varrho$ . By  $H^{-1}(D)$  we denote the dual of  $H_0^1(D)$ , the space of square integrable functions in  $D$  having zero trace on the boundary and square integrable distributional derivatives. We will often use the abbreviation  $\mathcal{H}_{\varrho}^2$  to denote the space

$$L^2_{\varrho}(\Gamma; H_0^1(D)) := \left\{ u : \bar{D} \times \Gamma \rightarrow \mathbb{R} : u \text{ strongly measurable and } \int_{\Gamma} \|u\|_{H_0^1(D)}^2 \varrho(\mathbf{y}) d\mathbf{y} < \infty \right\},$$

and  $\mathcal{H}_{\varrho}^{\infty}$  to denote the space

$$L^{\infty}_{\varrho}(\Gamma; H_0^1(D)) := \left\{ u : \bar{D} \times \Gamma \rightarrow \mathbb{R} : u \text{ strongly measurable and } \operatorname{ess\,sup}_{\mathbf{y} \in \Gamma} \|u(\cdot, \mathbf{y})\|_{H_0^1(D)} < \infty \right\}.$$

For the space  $H_0^1(D)$  we have the energy norm  $\|v\|_{H_0^1(D)} = \|\nabla v\|_{L^2(D)}$ , hence  $\mathcal{H}_{\varrho}^2$  is a Hilbert space with norm  $\|v\|_{\mathcal{H}_{\varrho}^2}^2 = \int_{\Gamma} \|v\|_{H_0^1(D)}^2 \varrho d\mathbf{y}$ . The *stochastic weak form* of problem (1) is given by: find  $u \in \mathcal{H}_{\varrho}^2$  such that  $\forall v \in \mathcal{H}_{\varrho}^2$

$$\int_{\Gamma} \mathcal{B}[u, v](\mathbf{y}) \varrho(\mathbf{y}) d\mathbf{y} = \int_{\Gamma} F(v) \varrho(\mathbf{y}) d\mathbf{y}, \quad (5)$$

where

$$\mathcal{B}[u, v](\mathbf{y}) = \int_D a(x, \mathbf{y}) \nabla u(x, \mathbf{y}) \cdot \nabla v(x, \mathbf{y}) dx, \quad F(v) = \int_D f(x) v(x, \mathbf{y}) dx. \quad (6)$$

For convenience, we will often use the abbreviation  $\mathcal{B}(\mathbf{y}) = \mathcal{B}[\cdot, \cdot](\mathbf{y})$  and suppress the dependence on  $x \in D$  in writing  $a(\mathbf{y}) = a(\cdot, \mathbf{y})$  and  $u(\mathbf{y}) = u(\cdot, \mathbf{y})$ . It follows from (A1) that  $\mathcal{B}(\mathbf{y})$  is a symmetric, uniformly coercive, and continuous bilinear operator on  $H_0^1(D)$ , parameterized by  $\mathbf{y} \in \Gamma$ , and  $\mathcal{B}(\mathbf{y})$  induces the norm

$$\|u\|_{\mathcal{B}(\mathbf{y})}^2 := \int_D a(x, \mathbf{y}) |\nabla u|^2 dx. \quad (7)$$

Assumption (A1) and the Lax-Milgram lemma also ensure the existence and uniqueness of the solution  $u$  to (5) in  $\mathcal{H}_{\varrho}^2$ .

The convergence of the global stochastic polynomial methods used to approximate (1) exploits the uniform ellipticity of the coefficient  $a(\mathbf{y})$  and depends on the regularity of  $u(\mathbf{y})$  with respect to  $\mathbf{y}$ . By  $\operatorname{Re}(z)$  and  $\operatorname{Im}(z)$  we denote the real and imaginary parts of  $z \in \mathbb{C}$ , and for  $0 < \delta < a_{\min}$  we define

$$U(a, \delta) = \{z \in \mathbb{C}^N : \operatorname{Re}(a(x, z)) \geq \delta, \forall x \in \bar{D}\}. \quad (8)$$

If  $U(a, \delta) \neq \emptyset$  for some  $0 < \delta < a_{\min}$ , we say that  $a(x, \mathbf{z})$  is uniformly elliptic on the set  $U(a, \delta)$  and we refer to  $U(a, \delta)$  as its domain of uniform ellipticity. For  $\boldsymbol{\gamma} = (\gamma_1, \dots, \gamma_N)$  with  $\gamma_i > 1 \ \forall i$  we denote the polyellipse

$$\mathcal{E}_{\boldsymbol{\gamma}} = \bigotimes_{1 \leq i \leq N} \left\{ z_i \in \mathbb{C} : \operatorname{Re}(z_i) \leq \frac{\gamma_i + \gamma_i^{-1}}{2} \cos \phi, \operatorname{Im}(z_i) \leq \frac{\gamma_i - \gamma_i^{-1}}{2} \sin \phi, \phi \in [0, 2\pi) \right\}.$$

In [33] it was shown that if  $a(\mathbf{y})$  satisfies (A1) and (A2), then for any  $0 < \delta < a_{\min}$  there exists a  $\boldsymbol{\gamma} = (\gamma_1, \dots, \gamma_N)$  with  $\gamma_i > 1 \ \forall i$  such that  $\mathcal{E}_{\boldsymbol{\gamma}} \subset U(a, \delta)$ . We can also similarly define the polydisc  $\mathcal{D}_{\boldsymbol{\gamma}} = \bigotimes_{1 \leq i \leq N} \{z_i \in \mathbb{C} : |z_i| \leq \gamma_i\}$ , though, for arbitrary  $0 < \delta < a_{\min}$ , it is not always possible to find a  $\boldsymbol{\gamma}$  with  $\gamma_i > 1 \ \forall i$  such that  $\mathcal{D}_{\boldsymbol{\gamma}} \subset U(a, \delta)$ . Figure 1 provides an illustration of this fact for various one-dimensional coefficients  $a(\mathbf{y})$ ,  $\mathbf{y} \in \mathbb{C}$ . Note that in the case of the 6th degree polynomial and exponential random variables, no disc of radius  $\gamma > 1$  containing  $\Gamma = [-1, 1]$  can fit in the region. The following theorem, proved in [33], shows the regularity of the solution  $u$  with respect to the parameterization.

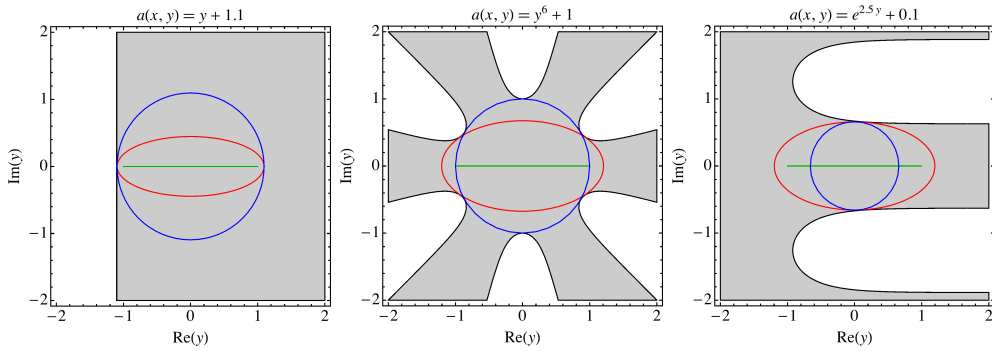


Figure 1: Domains of uniform ellipticity for some one-dimensional coefficients  $a(x, y)$  are indicated by the gray regions in each plot. The blue and red curves represent the maximal discs and ellipses, respectively, that can be contained in those domains, and the green lines represent the interval  $\Gamma = [-1, 1]$ .

**Theorem 2.4.** *When the coefficient  $a(x, \mathbf{y})$  satisfies (A1) and (A2), so that for some  $0 < \delta < a_{\min}$  and  $\boldsymbol{\gamma} = (\gamma_1, \dots, \gamma_N)$  with  $\gamma_i > 1 \ \forall i$  we have  $\mathcal{E}_{\boldsymbol{\gamma}} \subset U(a, \delta)$ , then the function  $\mathbf{z} \mapsto u(\mathbf{z})$  from (1) is holomorphic in an open neighborhood of  $\mathcal{E}_{\boldsymbol{\gamma}}$ .*

This result states that a direct consequence of the uniform ellipticity of the function  $a(x, \mathbf{y})$  on the polyellipse  $\mathcal{E}_{\boldsymbol{\gamma}} \subset U(a, \delta)$  is that the solution  $u$  of (1) has analytically smooth dependence on the parameterization  $\mathbf{y}$ . Theorem 2.4 is the key in motivating the construction of global stochastic Galerkin (SG) approximations to the solution  $u$  of (1), to be described in the following sections.

### 3. Stochastic Galerkin finite element method

In this section we define the SGFEM for constructing fully discrete approximations to the solution  $u$  of problem (1). This discretization employs mixed Galerkin projections in the spatial and parameter domains. In particular we rely on the finite element method for the spatial discretization, described in §3.1, and the stochastic Galerkin method for the parameter discretization, described in §3.2. In §3.3 we describe the linear systems that result from the SG discretization when Examples 2.1, 2.2, and 2.3 are used in problem (1). We then conclude in §3.4 with a discussion of the cost of solving the SG systems.

#### 3.1. Parameterized finite element approximation

We briefly define the finite element method for obtaining a discretization of  $u$  from (1) over the spatial domain  $D$ . Let  $\mathcal{T}_h$ , be a triangulation of  $D$  with maximum mesh size  $h > 0$ , and  $V_h(D) \subset H_0^1(D)$  a finite

element space of piecewise continuous polynomials on  $\mathcal{T}_h$  parameterized by  $h \rightarrow 0$ . Let  $\{\phi_j(x)\}_{j=1}^{J_h}$  denote a finite basis of  $V_h(D)$  of dimension  $J_h$ . We can write the *semi-discrete* problem as: find  $u_h(\mathbf{y}) \in V_h(D)$  such that  $\forall v \in V_h(D)$

$$\mathcal{B}[u_h(\mathbf{y}), v](\mathbf{y}) = F(v), \quad (9)$$

where  $\mathcal{B}[\cdot, \cdot](\mathbf{y})$  and  $F(\cdot)$  are defined in (6). For almost every  $\mathbf{y} \in \Gamma$ , problem (9) admits a unique solution of the form  $u_h(x, \mathbf{y}) = \sum_{j=1}^{J_h} u_j(\mathbf{y})\phi_j(x)$ . We discretize problem (9) by defining, for  $i, j = 1, \dots, J_h$ ,

$$[\mathbf{A}]_{i,j}(\mathbf{y}) = \mathcal{B}[\phi_j, \phi_i](\mathbf{y}), \quad \mathbf{F}_i = F(\phi_i). \quad (10)$$

The coefficients  $\mathbf{u}_h(\mathbf{y}) = [u_1(\mathbf{y}), u_2(\mathbf{y}), \dots, u_{J_h}(\mathbf{y})]^T$  of  $u_h(x, \mathbf{y})$  are determined by solving the linear system

$$\mathbf{A}(\mathbf{y})\mathbf{u}_h(\mathbf{y}) = \mathbf{F}, \quad (11)$$

at fixed realizations of  $\mathbf{y} \in \Gamma$ . Here  $\mathbf{A}(\mathbf{y})$  is symmetric and positive-definite so that (11) can be solved by iterative methods such as the conjugate gradient (CG) method.

### 3.2. Stochastic Galerkin approximation with an orthogonal basis

Based on the smoothness of the solution  $u$  to (1), characterized by Theorem 2.4, we now consider the construction of approximations to  $u$  in terms of global polynomials. Let  $\Lambda_p \subset \mathbb{N}_0^N$  be a finite set of multi-indices, e.g., having dimension  $\#\Lambda_p < \infty$ , and define the space of polynomials  $\mathcal{P}_{\Lambda_p}(\Gamma) = \text{span}\{\mathbf{y}^\nu : \nu \in \Lambda_p\}$ . A general global polynomial approximation problem can be framed in terms of solving for the  $\#\Lambda_p$  stochastic degrees of freedom (SDOF)  $\{u_{\mathbf{p}}\}_{\mathbf{p} \in \Lambda_p}$ . When an interpolatory approach is used, the resulting systems of equations are decoupled finite element systems. When Galerkin projection with an orthogonal basis is used, the finite element systems are fully coupled, and must be solved simultaneously. Some isotropic examples of such index sets include

$$\begin{aligned} \Lambda_p^{\text{TP}} &= \left\{ \mathbf{p} \in \mathbb{N}_0^N : \max_{1 \leq i \leq N} p_i \leq p \right\}, & \Lambda_p^{\text{TD}} &= \left\{ \mathbf{p} \in \mathbb{N}_0^N : \sum_{n=1}^N p_n \leq p \right\}, \\ \Lambda_p^{\text{SM}} &= \left\{ \mathbf{p} \in \mathbb{N}_0^N : \sum_{n=1}^N f(p_n) \leq f(p) \right\}, & f(p) &= \begin{cases} 0, & p = 0 \\ 1, & p = 1 \\ \lceil \log_2(p) \rceil, & p \geq 2 \end{cases} \end{aligned} \quad (12)$$

corresponding to the Tensor Products (TP), Total Degree (TD), and Smolyak (SM) polynomial spaces  $\mathcal{P}_{\Lambda_p^{\text{TP}}}(\Gamma)$ ,  $\mathcal{P}_{\Lambda_p^{\text{TD}}}(\Gamma)$ , and  $\mathcal{P}_{\Lambda_p^{\text{SM}}}(\Gamma)$ , respectively. When the solution  $u$  exhibits an anisotropic dependence on the parameters  $\mathbf{y}$ , anisotropic weighted versions of the index sets defined in (12) can be introduced to further reduce the number of SDOF needed to approximate  $u$  at a desired accuracy [3, 26].

**Remark 3.1** (Best  $M$ -term and quasi-optimal approximations). *The optimal choice of  $\Lambda_p$  would be the set  $\Lambda$  of cardinality  $M$  such that the corresponding approximation provides maximum accuracy out of all sets of size  $M$ . Such approximations are referred to as best  $M$ -term approximations, and recent work has focussed on the construction of best  $M$ -term Taylor and Galerkin approximations [4, 6, 7, 8, 33]. These approaches construct  $\Lambda$  by utilizing the largest  $M$  coefficients  $u_{\mathbf{p}}$  or sharp upper bounds of  $u_{\mathbf{p}}$ . However, in this effort we focus on analyzing the computational complexity of finding solutions to (5) in  $\mathcal{P}_{\Lambda_p}(\Gamma)$  for a prescribed index set  $\Lambda_p$ .*

For each  $n = 1, \dots, N$ , let  $\{\psi_{p_n}(y_n)\}_{p_n=1}^\infty$  be a sequence of univariate polynomials over  $\Gamma_n$ , orthonormal with respect to the  $L_{\varrho_n}^2(\Gamma_n)$  inner product. Then  $\{\Psi_{\mathbf{p}}(\mathbf{y})\}_{0 \leq |\mathbf{p}|}$  with  $\Psi_{\mathbf{p}}(\mathbf{y}) := \prod_{n=1}^N \psi_{p_n}(y_n)$  is a sequence of multivariate polynomials over  $\Gamma$ , orthonormal with respect to the  $L_{\varrho}^2(\Gamma)$  inner product. In the case that  $\varrho = \frac{1}{2}$  for each  $n = 1, \dots, N$ ,  $\{\psi_{p_n}\}_{p_n=1}^\infty$  and  $\{\Psi_{\mathbf{p}}\}_{0 \leq |\mathbf{p}|}$  are the univariate and multivariate Legendre polynomials, respectively. Given a specific choice of index set  $\Lambda_p$ , it follows that  $\{\Psi_{\mathbf{p}}\}_{\mathbf{p} \in \Lambda_p}$  forms a basis of

$\mathcal{P}_{\Lambda_p}(\Gamma)$  with dimension  $M_p = \dim(\mathcal{P}_{\Lambda_p}(\Gamma)) = \#\Lambda_p$ . Hence, with  $\{\phi_j\}_{j=1}^{J_h}$  as in §3.1 and  $\{\Psi_{\mathbf{p}}\}_{\mathbf{p} \in \Lambda_p}$  as above, we can now write the fully discrete stochastic Galerkin (SG) approximation as

$$u_{h,p}(x, \mathbf{y}) = \sum_{\mathbf{p} \in \Lambda_p} \sum_{j=1}^{J_h} u_{j,\mathbf{p}} \phi_j(x) \Psi_{\mathbf{p}}(\mathbf{y}), \quad (13)$$

whose coefficients can be found by solving the following coupled problem: *find  $u_{h,p} \in V_h(D) \otimes \mathcal{P}_{\Lambda_p}(\Gamma)$  such that for all  $v \in V_h(D) \otimes \mathcal{P}_{\Lambda_p}(\Gamma)$*

$$\mathbb{E} [\mathcal{B}[u_{h,p}, v](\mathbf{y})] = \mathbb{E} [F(v)], \quad (14)$$

where  $\mathcal{B}[\cdot, \cdot](\mathbf{y})$  and  $F(\cdot)$  are defined in (6). To form the linear system of equations resulting from the SG approximation given by (13), we let  $\mathbf{u}_{h,\mathbf{p}} = [u_{1,\mathbf{p}}, \dots, u_{J_h,\mathbf{p}}]^T$  be the vector of nodal values of the finite element solution corresponding to the  $\mathbf{p}$ -th stochastic mode of  $u_{h,p}$ , and  $\mathbf{u}_{h,p} = [\mathbf{u}_{h,\mathbf{p}}]_{\mathbf{p} \in \Lambda_p}^T$ . Observe that when  $f$  is deterministic  $\langle \Psi_{\mathbf{p}} \mathbf{F}_i \rangle = \mathbf{F}_i \delta_{\mathbf{0},\mathbf{p}}$  for all  $i = 1, \dots, J_h$ , where  $\delta_{\mathbf{0},\mathbf{p}} = 1$  if  $\mathbf{p} = \mathbf{0}$  and  $\delta_{\mathbf{0},\mathbf{p}} = 0$  otherwise. Performing a Galerkin projection onto  $\text{span}\{\Psi_{\mathbf{p}}\}_{\mathbf{p} \in \Lambda_p}$  for the solution of (14) yields the following system: for each  $\mathbf{p} \in \Lambda_p$

$$\sum_{\mathbf{q} \in \Lambda_p} \langle \Psi_{\mathbf{p}}(\mathbf{y}), \mathbf{A}(\mathbf{y}) \Psi_{\mathbf{q}}(\mathbf{y}) \rangle \mathbf{u}_{h,\mathbf{q}} = \langle \Psi_{\mathbf{p}}(\mathbf{y}), \mathbf{F} \rangle, \quad (15)$$

which can be written algebraically as a system of *fully coupled finite element problems*: for each  $\mathbf{p} \in \Lambda_p$

$$\sum_{\mathbf{q} \in \Lambda_p} [\mathbf{K}]_{\mathbf{p},\mathbf{q}} \mathbf{u}_{h,\mathbf{q}} = \mathbf{F} \delta_{\mathbf{0},\mathbf{p}} \quad (16)$$

with  $[\mathbf{K}]_{\mathbf{p},\mathbf{q}} = \langle \Psi_{\mathbf{p}}(\mathbf{y}), \mathbf{A}(\mathbf{y}) \Psi_{\mathbf{q}}(\mathbf{y}) \rangle$  and  $\mathbf{A}(\mathbf{y})$  as given in (10).

**Remark 3.2.** Typically matrix free methods are applied to solve (16) without ever explicitly forming  $\mathbf{K}$  in memory, as described in [27]. When the resulting system is sparse, as a result of an affine coefficient  $a(x, \mathbf{y})$ , e.g., Example 2.1, this can lead to computationally efficient solution strategies. However, these implementations rely on the fact that the coefficient  $a(x, \mathbf{y})$  can be written as a sum of separable functions of  $x$  and  $\mathbf{y}$ , e.g.,  $a(x, \mathbf{y}) = \sum_{j=1}^N b_j(x) c_j(\mathbf{y})$ . For the transcendental function  $a(x, \mathbf{y})$  from Example 2.3, this may not be the case. Moreover, when  $\mathbf{K}$  is block-dense, matrix-vector multiplications require approximately  $\mathcal{O}(J_h M_p^2)$  floating point operations (FLOPs), so that when iterative methods are used, the solution of the fully coupled finite element problems given in (16) becomes unfeasible.

### 3.3. Representations of $a(\mathbf{y})$ and the corresponding matrix $\mathbf{K}$

For a general coefficient  $a(x, \mathbf{y})$ , the matrix  $\mathbf{K}$  in (16) requires the storage of at most  $M_p^2$  block matrices of the size and sparsity of  $\mathbf{A}(\mathbf{y})$ , i.e.,  $\mathcal{O}(J_h M_p^2)$  elements. However, in several specific cases the actual block-sparsity of  $\mathbf{K}$  is much less. We recall the coefficient from Example 2.1, where  $\mathbf{K}$  can be rewritten

$$[\mathbf{K}]_{\mathbf{p},\mathbf{q}} = \langle \Psi_{\mathbf{p}}(\mathbf{y}), \Psi_{\mathbf{q}}(\mathbf{y}) \rangle \mathbf{A}_0 + \sum_{k=1}^N \langle y_k \Psi_{\mathbf{p}}(\mathbf{y}), \Psi_{\mathbf{q}}(\mathbf{y}) \rangle \mathbf{A}_k,$$

with  $[\mathbf{A}_0]_{i,j} = \int_D a_0(x) \nabla \phi_j(x) \cdot \nabla \phi_i(x) dx$  and  $[\mathbf{A}_k]_{i,j} = \int_D b_k(x) \nabla \phi_j(x) \cdot \nabla \phi_i(x) dx$ . If we let  $[\mathbf{G}_0]_{\mathbf{p},\mathbf{q}} = \langle \Psi_{\mathbf{p}}(\mathbf{y}), \Psi_{\mathbf{q}}(\mathbf{y}) \rangle$  and  $[\mathbf{G}_k]_{\mathbf{p},\mathbf{q}} = \langle y_k \Psi_{\mathbf{p}}(\mathbf{y}), \Psi_{\mathbf{q}}(\mathbf{y}) \rangle$ , then  $\mathbf{K}$  has a matrix representation, given by,

$$\mathbf{K} = \mathbf{G}_0 \otimes \mathbf{A}_0 + \sum_{k=1}^N \mathbf{G}_k \otimes \mathbf{A}_k, \quad (17)$$

where  $\mathbf{A} \otimes \mathbf{B}$  denotes the Kronecker product of  $\mathbf{A}$  and  $\mathbf{B}$ . We note that a similar construction can be obtained for any coefficient  $a(x, \mathbf{y})$  which can be written as a sum of separable functions of  $x$  and  $\mathbf{y}$ , such as the polynomial function of Example 2.2.

However, when  $a(x, \mathbf{y})$  is not separable in  $x$  and  $\mathbf{y}$ , this construction is no longer valid, and the resulting matrix  $\mathbf{K}$  may be block-dense if we simply carry out the Galerkin projections and compute  $\mathbf{K}$  directly. For certain special cases, e.g., when the diffusion coefficient is given by a log-transformed random field, the problem can be reformulated as a convection-diffusion problem, and the resulting system can be solved much more efficiently than the original problem [35]. In general, this reformulation is not applicable. Hence, for a general transcendental coefficient  $a(x, \mathbf{y})$ , as given in Example 2.3, we project the coefficient onto an additional subspace  $\mathcal{P}_{\Lambda_r}(\Gamma)$ ,  $r \in \mathbb{N}_0$ , in order to obtain a separable representation. To see this, define  $\{\Psi_r(\mathbf{y})\}_{0 \leq |r|}$  to be the (infinite) basis of orthonormal polynomials of  $L^2_\rho(\Gamma)$  as in §3.2. Then  $a(x, \mathbf{y})$  can be written as an expansion such that  $a(x, \mathbf{y}) = \sum_{0 \leq |r|} a_r(x) \Psi_r(\mathbf{y})$ , where the coefficients  $a_r(x) = \langle a(x, \mathbf{y}), \Psi_r(\mathbf{y}) \rangle$ . Let  $\Lambda_r$  be an index set of the type described in §3.2. Since infinite series representations are not practical in computations, we seek a truncation

$$a^r(x, \mathbf{y}) := \sum_{\mathbf{r} \in \Lambda_r} a_r(x) \Psi_{\mathbf{r}}(\mathbf{y}) \quad (18)$$

in the subspace  $\mathcal{P}_{\Lambda_r}(\Gamma)$  for some  $r \in \mathbb{N}_0$ . When  $a^r(x, \mathbf{y}) \neq a(x, \mathbf{y})$ , e.g., in the case that the projection order  $r$  is chosen to minimize error independent of the SG discretization, we let  $u_{h,p}^r$  denote the corresponding solution to the fully discrete SG approximation problem with  $a(x, \mathbf{y})$  replaced with  $a^r(x, \mathbf{y})$ . By substituting  $a^r(x, \mathbf{y})$  into (6) we obtain

$$\int_D \left( \sum_{\mathbf{r} \in \Lambda_r} a_r(x) \Psi_{\mathbf{r}}(\mathbf{y}) \right) \nabla \phi_j(x) \cdot \nabla \phi_i(x) dx = \sum_{\mathbf{r} \in \Lambda_r} [\mathbf{A}_r]_{i,j} \Psi_{\mathbf{r}}(\mathbf{y}), \quad (19)$$

$$[\mathbf{A}_r]_{i,j} = \int_D a_r(x) \nabla \phi_j(x) \cdot \nabla \phi_i(x) dx. \quad (20)$$

Equation (19) represents an expansion of the stochastic finite element stiffness matrix  $\mathbf{A}(\mathbf{y})$  and equation (20) represents the  $\mathbf{r}$ -th mode of the expansion. Let  $\mathbf{u}_{h,p}^r = [u_{1,p}^r, \dots, u_{j_h,p}^r]^T$  denote the vector of nodal values of the finite element solution corresponding to the  $\mathbf{p}$ -th stochastic mode of  $u_{h,p}^r$ , and  $\mathbf{u}_{h,p}^r = [\mathbf{u}_{h,p}^r]_{\mathbf{p} \in \Lambda_p}^T$ . We substitute the expansion of  $\mathbf{A}(\mathbf{y})$  into the Galerkin equations (15), to obtain the coupled system: for each  $\mathbf{p} \in \Lambda_p$

$$\sum_{\mathbf{r} \in \Lambda_r} \sum_{\mathbf{q} \in \Lambda_p} [\mathbf{G}_r]_{\mathbf{p},\mathbf{q}} \mathbf{A}_r \mathbf{u}_{h,q}^r = \langle \Psi_{\mathbf{p}}, \mathbf{F} \rangle, \quad [\mathbf{G}_r]_{\mathbf{p},\mathbf{q}} = \langle \Psi_{\mathbf{p}} \Psi_{\mathbf{q}} \Psi_{\mathbf{r}} \rangle. \quad (21)$$

Alternatively, similar to (17), we may define  $\mathbf{K}_r = \sum_{\mathbf{r} \in \Lambda_r} \mathbf{G}_r \otimes \mathbf{A}_r$  to again obtain the coupled system of finite element problems: for all  $\mathbf{p} \in \Lambda_p$

$$\sum_{\mathbf{q} \in \Lambda_p} [\mathbf{K}_r]_{\mathbf{p},\mathbf{q}} \mathbf{u}_{h,q}^r = \mathbf{F} \delta_{\mathbf{0},\mathbf{p}} \quad \forall \mathbf{p} \in \Lambda_p. \quad (22)$$

Note that in forming  $\mathbf{K}_r$ , we now need only store the matrices  $\{\mathbf{G}_r\}_{r \in \Lambda_r}$  and  $\{\mathbf{A}_r\}_{r \in \Lambda_r}$ , so that the efficient, matrix-free, solution strategies discussed in Remark 3.2 can be applied.

**Remark 3.3** (Projection and well-posedness). *In the case that the coefficient is a transcendental function of the random variables, as in Example 2.3, there does not exist a  $r \in \mathbb{N}_0$  such that the projection (18) is exact. Due to the orthogonality of the basis, setting  $r = 2p$  in the construction of  $\mathbf{K}_r$  yields an entirely block-dense system [23] that is equivalent to (16), and computationally infeasible to solve. A more practical approach is to choose the expansion order  $0 \leq r \leq 2p$ , based on a-priori estimates of the error in the solution introduced by the truncation, so that the error when using the truncated expansion does not exceed that of*



the SG approximation. In this approach however, it becomes important to consider whether the truncated projection violates the well-posedness of (1) by failing to satisfy assumption (A1). One way to guarantee this is to choose  $\tilde{r} \leq r \leq 2p$  such that

$$\tilde{r} := \min\{r \in \mathbb{N}_0 : \|a - a^r\|_{L^\infty_\varphi(\Gamma; L^\infty(D))} \leq a_{\min}, \forall \nu \in \mathbb{N}_0, \nu \geq r\}. \quad (23)$$

An example of this problem can be seen in Figure 2 where for the function  $a(x, y) = 0.1 + \exp(2.5y)$ , uniform ellipticity of the truncated projection  $a^r(x, y)$  does not hold on  $\Gamma = [-1, 1]$  for  $r < 4$ .

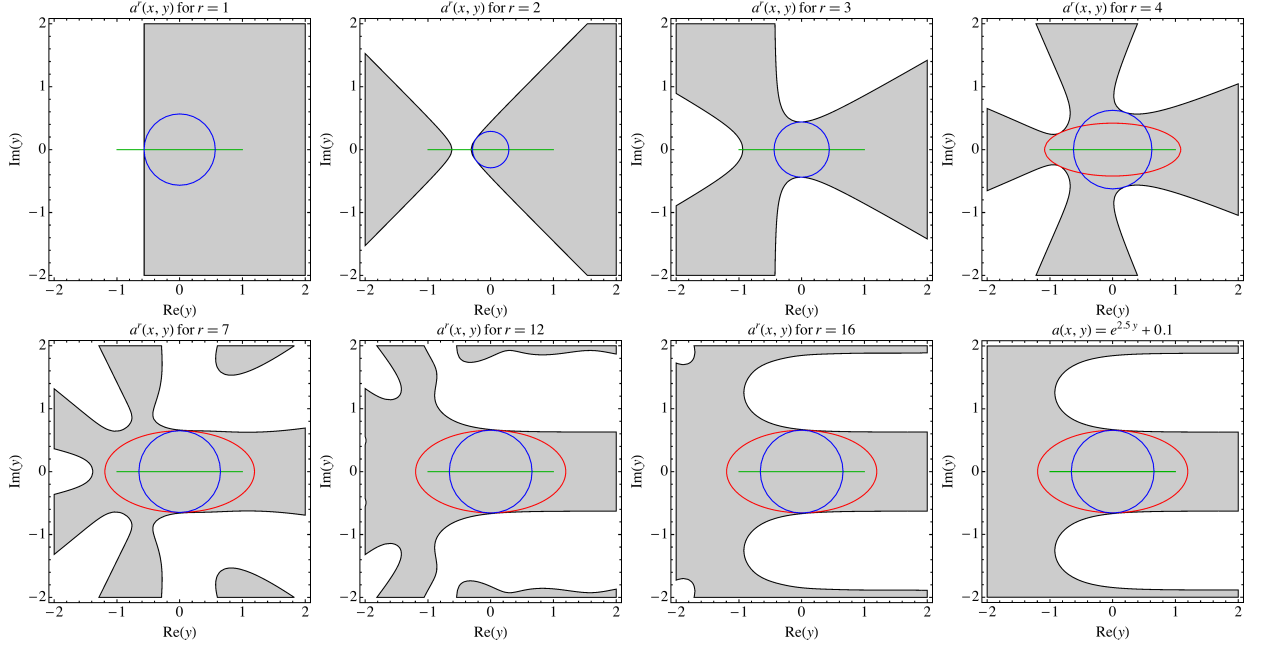


Figure 2: Domains of uniform ellipticity for the total degree orthogonal expansions of order  $r$  of the one-dimensional coefficient  $a(x, y) = 0.1 + \exp(2.5y)$ , for  $y \in \Gamma = [-1, 1] \subset \mathbb{R}^1$ , are indicated by the gray regions in each plot. The last plot shows the domain of uniform ellipticity of the original function  $a(x, y)$ . The blue and red curves represent the maximal discs and ellipses, respectively, that can be contained in those domains, and the green lines represent the interval  $\Gamma$ .

### 3.4. Cost of solving the generalized SG system

Without loss of generality, to solve the stochastic Galerkin system (22) for  $\mathbf{u}_{h,p}^r$ , we use the preconditioned conjugate gradient (PCG) method, wherein, for the unpreconditioned CG method, we have the estimate

$$\|\mathbf{u}_{h,p}^r - \mathbf{u}_{h,p}^{r,(k)}\|_{\mathbf{K}_r} \leq 2 \left( \frac{\sqrt{\kappa_r} - 1}{\sqrt{\kappa_r} + 1} \right)^k \|\mathbf{u}_{h,p}^r - \mathbf{u}_{h,p}^{r,(0)}\|_{\mathbf{K}_r}. \quad (24)$$

Here  $\kappa_r$  is the condition number of  $\mathbf{K}_r$ ,  $\mathbf{u}_{h,p}^{r,(0)}$  is the vector of the initial guess, and  $\mathbf{u}_{h,p}^{r,(k)}$  is the output of the  $k$ -th iteration of the CG solver. The CG method is highly dependent on the conditioning of the system, and when  $\kappa_r$  is large, the number of iterations needed to reduce the error in  $\mathbf{u}_{h,p}^{r,(k)}$  will also be significant. Hence we introduce the mean-based block-diagonal preconditioner (see, e.g., [27, 29]),

$$\mathbf{P} := \mathbf{G}_0 \otimes \mathbf{A}_0, \quad (25)$$

with  $\mathbf{A}_0$  and  $\mathbf{G}_0$  the matrices defined in (20) and (21) for  $\mathbf{r} = \mathbf{0}$ , respectively.

For  $\mathbf{r} \in \Lambda_r$ , at every iteration of the CG method, or any iterative approach, each nonzero entry in each matrix  $\mathbf{G}_{\mathbf{r}}$  implies a matrix-vector product of the form  $\langle \Psi_{\mathbf{p}} \Psi_{\mathbf{q}} \Psi_{\mathbf{r}} \rangle \mathbf{A}_{\mathbf{r}} \mathbf{p}_{\mathbf{q}}^{(k)}$ , where  $\langle \Psi_{\mathbf{p}} \Psi_{\mathbf{q}} \Psi_{\mathbf{r}} \rangle$  is a scalar quantity. Let  $\text{nnz}(\mathbf{A})$  denote the number of nonzeros of a matrix  $\mathbf{A}$ , and define

$$\mathcal{M}(p, r) = \sum_{\mathbf{r} \in \Lambda_r} \text{nnz}(\mathbf{G}_{\mathbf{r}}) \quad (26)$$

to be the total number of nonzeros in all of the matrices  $\{\mathbf{G}_{\mathbf{r}}\}_{\mathbf{r} \in \Lambda_r}$  at order  $p$ . With this in mind, an upper bound for the work in floating point operations (FLOPs) of solving (22) is given by

$$W^{\text{SG}} \approx \mathcal{O}(J_h) * \mathcal{M}(p, r) * N_{\text{iter}}^{\text{SG}}, \quad (27)$$

where the term  $\mathcal{O}(J_h)$  corresponds to the cost of a single finite element matrix-vector product, and  $N_{\text{iter}}^{\text{SG}}$  is the number of iterations of the CG solver without a preconditioner. If we apply a preconditioner, in hopes to minimize  $N_{\text{iter}}^{\text{SG}}$ , we must also account for the added cost of applying the preconditioner at each iteration. With the mean-based preconditioner from (25), at each iteration we multiply an additional matrix of size  $J_h M_p \times J_h M_p$ , but the matrix consists only of  $M_p$  diagonal blocks. Here we assume that in finding the inverse of the preconditioning matrix  $\mathbf{P}$ , a sparse approximation to  $\mathbf{P}^{-1}$  is used, which we will denote  $\tilde{\mathbf{P}}^{-1}$ . Such a decomposition can be found from, e.g., incomplete LU or incomplete Cholesky factorizations. Hence, for each iteration we require  $M_p$  additional matrix-vector products of the size, and complexity, of the original finite element system, so the work estimate in FLOPs for the case of this preconditioner is given by

$$W^{\text{pSG}} \approx \mathcal{O}(J_h) * (M_p + \mathcal{M}(p, r)) * N_{\text{iter}}^{\text{pSG}}, \quad (28)$$

where  $N_{\text{iter}}^{\text{pSG}}$  is the number of iterations needed by the PCG method. Other preconditioners, such as the Kronecker product preconditioner suggested in [34] would require a different form of (28).

Figure 3 displays the effect of fixing the projection order of the solution but increasing the order of the projection of the coefficient. In order to minimize the error of the projection, such a situation would be required if the coefficient is highly nonlinear and reflects the importance of considering  $\mathcal{M}(p, r)$  in the computational cost of the SGFEM.

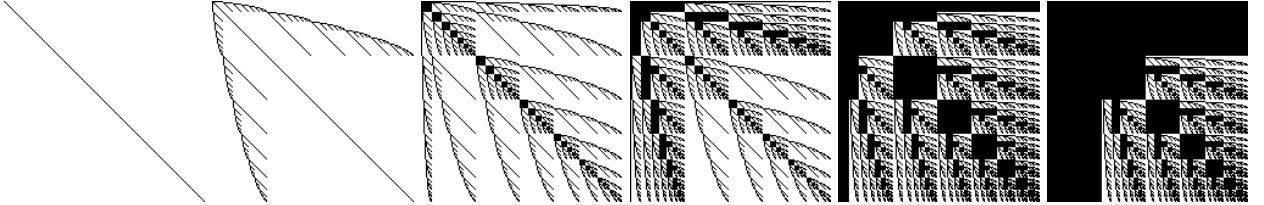


Figure 3: Visualization of the number of nonzeros of a  $165 \times 165$  SG matrix with elements  $[\mathbf{K}_r]_{\mathbf{p}, \mathbf{q}} = \sum_{\mathbf{r} \in \Lambda_r} [\mathbf{G}_{\mathbf{r}}]_{\mathbf{p}, \mathbf{q}} * \mathbf{A}_{\mathbf{r}}$ . Each pixel represents a block finite element system when using a total degree projection of the solution of fixed degree  $p = 3$ , and increasing the total degree of the projection of the coefficient, i.e.,  $r = 0, 1, 2, 3, 4, 5$ . At  $r = 6$ , the matrix is entirely block-dense.

#### 4. Explicit cost bounds for the SGFEM

The primary goal of this section is to estimate the algorithmic complexity required by the SGFEM to construct an approximation to (1) within a prescribed tolerance  $\varepsilon > 0$ . We assume  $a(x, \mathbf{y})$  is a general non-affine coefficient, as in Examples 2.2 and 2.3, satisfying assumptions (A1) and (A2). Let  $a^r(x, \mathbf{y})$  be the orthogonal expansion of  $a(x, \mathbf{y})$ , given by (18), of total degree  $r$ , i.e.,  $a^r(x, \mathbf{y}) \in \mathcal{P}_{\Lambda_r}(\Gamma)$  with  $\Lambda_r = \Lambda_r^{\text{TD}}$  from (12). We further assume that  $\tilde{r} \leq r \leq 2p$ , with  $\tilde{r}$  given in (23), so that  $a^r(x, \mathbf{y})$  also satisfies (A1) and (A2). We will focus on the complexity of solving (22), when the stochastic discretization to (1) is performed in

$\mathcal{P}_{\Lambda_p}(\Gamma)$  with  $\Lambda_p = \Lambda_p^{\text{TD}}$  from (12), i.e., in the space of total degree polynomials of order  $p$ , and the physical discretization is performed with the finite element method. These results are presented in the context of solving the linear system (22) with a PCG method when a zero initial vector is used to seed the solver. The results, however, can be generalized to other methods, such as preconditioned GMRES and other Krylov subspace methods.

The results are organized as follows. In §4.1 we discuss the overall complexity of the matrix-vector products associated with solving (22) when using the SG matrix  $\mathbf{K}_r = \sum_{\mathbf{r} \in \Lambda_r} \mathbf{G}_r \otimes \mathbf{A}_r$ . Our analysis extends the results of [14] in order to provide a bound on the block-sparsity of the SG system  $\mathbf{K}_r$  in the more general setting of a non-affine coefficient  $a^r(x, \mathbf{y})$ , as given in Examples 2.2 and 2.3. In particular, for  $\mathbf{G}_r$  given in (21), we show that  $\text{nnz}(\mathbf{G}_r) = \mathcal{O}(\min\{2^{|\mathbf{r}|}, M_{\lceil |\mathbf{r}|/2 \rceil}\} M_{p - \lceil |\mathbf{r}|/2 \rceil})$  for every  $\mathbf{r} \in \Lambda_r$ , where  $M_r = \binom{N+r}{N}$  for  $r \in \mathbb{N}_0$ , when solving (22), so that the total complexity of the matrix-vector products with the Galerkin system is  $\mathcal{O}(J_h M_p M_r \min\{2^r, M_{\lceil r/2 \rceil}\})$ . In §4.2, we perform an  $\varepsilon$ -complexity analysis to derive the explicit cost bounds of the SGFEM using PCG, in terms of FLOPs as the tolerance  $\varepsilon \rightarrow 0$ . Finally, in §4.3 we discuss issues related to the conditioning of the SG system.

#### 4.1. Complexity of matrix-vector multiplications in the SG approximation

In this section we provide rigorous bounds on the sparsity of the SG matrix  $\mathbf{K}_r$  from (22), for arbitrary  $0 \leq r \leq 2p$ , and  $p \in \mathbb{N}_0$ . Our main result, given by Theorem 4.1, provides an exact count for  $\text{nnz}(\mathbf{G}_r)$  in the general case  $|\mathbf{r}| \in \mathbb{N}_0$  and  $N \in \mathbb{N}$  when the integrals  $[\mathbf{G}_r]_{\mathbf{p}, \mathbf{q}} = \langle \Psi_r \Psi_p \Psi_q \rangle$  are defined in terms of even weight functions  $\varrho$ . This result can be seen as an extension of estimates from [14], where bounds on the sparsity of  $\mathbf{G}_r$  were shown in the cases (i)  $|\mathbf{r}| = 1$  and  $N \in \mathbb{N}$  and (ii)  $N = 1$  and  $\mathbf{r} = r \in \mathbb{N}_0$ . We then provide upper bounds on the total number of nonzeros blocks  $\mathcal{M}(p, r) = \sum_{\mathbf{r} \in \Lambda_r} \text{nnz}(\mathbf{G}_r)$  of the matrix  $\mathbf{K}_r$  from (22), both in the cases that  $a(x, \mathbf{y})$  is a finite order polynomial as in Example 2.2 and the case that  $a(x, \mathbf{y})$  is a transcendental function of the random variables, as in Example 2.3. Our first major result is summarized in the following Theorem:

**Theorem 4.1.** *Let  $\Lambda_p$  and  $\Lambda_r$  be the isotropic total degree index sets corresponding to the solution and coefficient, respectively, with  $p, r \in \mathbb{N}_0$ , and  $0 \leq r \leq 2p$ . If  $\mathbf{r} \in \Lambda_r$ , and  $\varrho_i$  are even for all  $i = 1, \dots, N$ , then for the matrix  $\mathbf{G}_r$  from (21) we have*

$$\text{nnz}(\mathbf{G}_r) = \sum_{\ell=\lceil |\mathbf{r}|/2 \rceil}^{|\mathbf{r}|} c(\mathbf{r}, \ell) \binom{N+p-\ell}{p-\ell}, \quad c(\mathbf{r}, \ell) = \begin{cases} \#\mathbf{S}(\mathbf{r}, \ell) & |\mathbf{r}| \text{ even}, \ell = |\mathbf{r}|/2, \\ 2\#\mathbf{S}(\mathbf{r}, \ell) & \text{otherwise}, \end{cases} \quad (29)$$

with  $\mathbf{S}(\mathbf{r}, \ell) = \{\mathbf{s} \in \mathbb{N}_0^N : |\mathbf{s}| = \ell, \mathbf{s} \leq \mathbf{r}\}$ , so that  $\#\mathbf{S}(\mathbf{r}, \ell)$  is equal to the coefficient of  $t^\ell$  in the polynomial  $P_{\mathbf{r}}(t) = \prod_{i=1}^N \sum_{j=0}^{r_i} t^j$ . Moreover, we have the following bound for  $\text{nnz}(\mathbf{G}_r)$ , i.e.,

$$\text{nnz}(\mathbf{G}_r) \leq 2 \min \left\{ 2^{|\mathbf{r}|}, \binom{N + \lceil |\mathbf{r}|/2 \rceil}{N} \right\} \binom{N+p - \lceil |\mathbf{r}|/2 \rceil}{N}, \quad (30)$$

so that

$$\mathcal{M}(p, r) \leq 2 \sum_{j=0}^r \min \left\{ 2^j, \binom{N + \lceil j/2 \rceil}{N} \right\} \binom{N-1+j}{N-1} \binom{N+p - \lceil j/2 \rceil}{N}. \quad (31)$$

PROOF. For a given  $\mathbf{r} \in \Lambda_r$ , we estimate the number of pairs  $(\mathbf{p}, \mathbf{q}) \in \Lambda_p \times \Lambda_p$  such that  $\langle \Psi_r \Psi_p \Psi_q \rangle = \prod_{i=1}^N \langle \psi_{r_i} \psi_{p_i} \psi_{q_i} \rangle \neq 0$ . To do this, we extend the result of [14, Lemma 28] to a general matrix  $\mathbf{G}_r$  with  $|\mathbf{r}| \in \mathbb{N}$ . Since  $\{\Psi_r\}_{\mathbf{r} \in \Lambda_r}$  are orthonormal with respect to the even weight function  $\rho(\mathbf{y}) = \prod_{i=1}^N \rho_i(y_i)$ , we see that  $\langle \Psi_r \Psi_p \Psi_q \rangle \neq 0$  if and only if  $(\mathbf{p}, \mathbf{q}) \in \Theta_r$ , where

$$\Theta_r = \{(\mathbf{p}, \mathbf{q}) \in \Lambda_p \times \Lambda_p : |p_i - q_i| \leq r_i \leq p_i + q_i, \\ \text{and } p_i + q_i + r_i \text{ is even } \forall i = 1, \dots, N\}.$$

Therefore, to estimate the number of nonzeros in the matrix  $\mathbf{G}_{\mathbf{r}}$ , we must estimate  $\#\Theta_{\mathbf{r}}$ . However,  $\Theta_{\mathbf{r}}$  is different for each  $\mathbf{r} \in \Lambda_r$ . Even when  $\mathbf{r}_1, \mathbf{r}_2 \in \Lambda_r$  are such that  $|\mathbf{r}_1| = |\mathbf{r}_2|$ , in general we do not have that  $\#\Theta_{\mathbf{r}_1} = \#\Theta_{\mathbf{r}_2}$ . On the other hand, if  $\mathbf{r}_2$  is a permutation of  $\mathbf{r}_1$ , then it is easy to see that  $\#\Theta_{\mathbf{r}_1} = \#\Theta_{\mathbf{r}_2}$  since  $\Lambda_p$  is the isotropic total degree set. Also note that  $\langle \Psi_{\mathbf{r}} \Psi_{\mathbf{p}} \Psi_{\mathbf{q}} \rangle = \langle \Psi_{\mathbf{r}} \Psi_{\mathbf{q}} \Psi_{\mathbf{p}} \rangle$  so if  $(\mathbf{p}, \mathbf{q}) \in \Theta_{\mathbf{r}}$  then  $(\mathbf{q}, \mathbf{p}) \in \Theta_{\mathbf{r}}$  as well. Note that  $\Theta_{\mathbf{r}}$  can be rewritten

$$\Theta_{\mathbf{r}} = \{(\mathbf{p}, \mathbf{q}) \in \Lambda_p \times \Lambda_p : |p_i - q_i| \leq r_i \leq p_i + q_i, \\ \text{and } |p_i - q_i| + r_i \text{ is even } \forall i = 1, \dots, N\}.$$

Hence if  $(\mathbf{p}, \mathbf{q}) \in \Theta_{\mathbf{r}}$ , then we see that  $(\mathbf{p}, \mathbf{q})$  must satisfy

- (i)  $|p_i - q_i| \leq r_i$  for all  $i = 1, \dots, N$ ,
- (ii)  $r_i \leq p_i + q_i$  for all  $i = 1, \dots, N$ , and
- (iii)  $|p_i - q_i| + r_i$  is even for all  $i = 1, \dots, N$ .

Note that when  $r_i = 0$ , we see that  $p_i = q_i \leq p$ , and when  $r_i > 0$  we see that (i) and (iii) imply  $|p_i - q_i| \in \{0, 2, 4, \dots, r_i\}$  for  $r_i$  even, and  $|p_i - q_i| \in \{1, 3, 5, \dots, r_i\}$  for  $r_i$  odd. For each  $i = 1, \dots, N$ , let  $\{k_i^{(n)}\}_{n=0}^{\lfloor r_i/2 \rfloor}$  be the sequence defined by

$$k_i^{(n)} = \begin{cases} 2n + 1 & r_i \text{ odd,} \\ 2n & r_i \text{ even,} \end{cases}$$

so that fixing  $|p_i - q_i| = k_i^{(n)}$  implies that conditions (i) and (iii) are met.

To satisfy (ii) we must have  $r_i \leq p_i + q_i$  and to satisfy (i) and (iii) we must have  $|p_i - q_i| = k_i^{(n)}$ . To avoid overcounting due to symmetry, we first fix possible values of  $p_i$  and consider what  $q_i$  must be. Let  $\{s_i^{(n)}\}_{n=0}^{\lfloor r_i/2 \rfloor}$  be the sequence defined by

$$s_i^{(n)} = \frac{r_i + k_i^{(n)}}{2},$$

which we will refer to as the sequence of *starting points* for  $p_i$  corresponding to  $k_i^{(n)}$ . Note that the starting points  $\{s_i^{(n)}\}_{n=0}^{\lfloor r_i/2 \rfloor}$  enumerate the integers between  $\lceil r_i/2 \rceil$  and  $r_i$ . Picking  $p_i \in \{s_i^{(n)}, s_i^{(n)} + 1, \dots, p\}$  and  $q_i = p_i - k_i^{(n)}$  we have

$$p_i + q_i = 2p_i - k_i^{(n)} \geq 2s_i^{(n)} - k_i^{(n)} = 2 \left( \frac{r_i + k_i^{(n)}}{2} \right) - k_i^{(n)} = r_i$$

or  $p_i + q_i \geq r_i$ , so that (ii) is satisfied.

Since (i), (ii), and (iii) are satisfied by setting  $p_i \in \{s_i^{(n)}, s_i^{(n)} + 1, \dots, p\}$  and  $q_i = p_i - k_i^{(n)}$  for a fixed  $0 \leq n \leq \lfloor r_i/2 \rfloor$ , we count the number of admissible pairs for these choices. In  $N - 1$  variables, the number of polynomials of total degree less than or equal to  $p - p_i$  is given by

$$\binom{N - 1 + p - p_i}{p - p_i},$$

where  $\binom{n}{k} = 0$  if  $n < k$  or  $n, k < 0$ . To simplify notation, pick  $s_i = s_i^{(n)}$  (one of the starting points in the  $i$ -th direction) and  $k_i = k_i^{(n)}$  (its associated distance), where  $0 \leq n \leq \lfloor r_i/2 \rfloor$  is fixed. To count the number of admissible pairs associated with the difference  $k_i$  and starting point  $s_i$ , we compute

$$\sum_{p_i=s_i}^p \binom{N - 1 + p - p_i}{p - p_i} = \sum_{j=0}^{p-s_i} \binom{N - 1 + j}{j} = \binom{N + p - s_i}{p - s_i}.$$

Define  $\mathbf{s} \in \mathbb{N}_0^N$  with the  $s_i$  as above, then  $\mathbf{s}$  corresponds to a possible combination of starting points in each direction. To estimate the number of polynomials associated with the starting point  $\mathbf{s}$ , we compute

$$\sum_{p_1=s_1}^p \sum_{p_2=s_2}^{p-p_1} \cdots \sum_{p_N=s_N}^{p-p_1-\cdots-p_{N-1}} \binom{p-p_1-\cdots-p_N}{p-p_1-\cdots-p_N} = \binom{N+p-|\mathbf{s}|}{p-|\mathbf{s}|}, \quad (32)$$

where the sum easily follows by an induction argument and Pascal's rule.

Enumerating all of the pairs  $(\mathbf{p}, \mathbf{q}) \in \Theta_{\mathbf{r}}$  thus reduces to counting the number of possible combinations of starting points. Hence, in  $N$  dimensions we consider all such multi-indices of the  $\{s_i^{(n)}\}_{n=0}^{\lceil r_i/2 \rceil}$  whose components sum to some integer  $\lceil |\mathbf{r}|/2 \rceil \leq \ell \leq |\mathbf{r}|$ . For two multi-indices  $\mathbf{s}, \mathbf{r} \in \mathbb{N}_0^N$ , we say  $\mathbf{s} \leq \mathbf{r}$  if and only if  $s_i \leq r_i$  for all  $i = 1, \dots, N$ . Define the set  $\mathbf{S}(\mathbf{r}, \ell) = \{\mathbf{s} \in \mathbb{N}_0^N : |\mathbf{s}| = \ell, \mathbf{s} \leq \mathbf{r}\}$ , which corresponds to a particular slice of the desired set of starting points. To estimate  $\#\mathbf{S}(\mathbf{r}, \ell)$ , we consider the familiar counting argument of placing  $N$  bars among  $\ell$  stars with the added restriction that the number of stars in the  $i$ -th bin not exceed  $r_i$ . Such a problem can be reframed in terms of finding the coefficient  $c(\mathbf{r}, \ell)$  of  $t^\ell$  in the generating function  $P_{\mathbf{r}}(t) = \prod_{i=1}^N \sum_{j=0}^{r_i} t^j$ . Combining (32) and summing over  $\ell$  between  $\lceil |\mathbf{r}|/2 \rceil \leq \ell \leq |\mathbf{r}|$  we arrive at (29), where the coefficients  $c(\mathbf{r}, \ell) = \#\mathbf{S}(\mathbf{r}, \ell)$  when  $|\mathbf{r}|$  is even and  $\ell = |\mathbf{r}|/2$  (in this case the roles of  $p_i$  and  $q_i$  can not be reversed) and  $c(\mathbf{r}, \ell) = 2\#\mathbf{S}(\mathbf{r}, \ell)$  otherwise.

Noting that  $\cup_{\ell=\lceil |\mathbf{r}|/2 \rceil}^{|\mathbf{r}|} \mathbf{S}(\mathbf{r}, \ell)$  is a change of coordinates of a total degree index set of order  $\lceil |\mathbf{r}|/2 \rceil$  intersected with the hyperrectangle  $\{\mathbf{s} \in \mathbb{N}_0^N : \mathbf{s} \leq \mathbf{r}\}$  yields the bound

$$\sum_{\ell=\lceil |\mathbf{r}|/2 \rceil}^{|\mathbf{r}|} c(\mathbf{r}, \ell) \leq 2 \sum_{\ell=\lceil |\mathbf{r}|/2 \rceil}^{|\mathbf{r}|} \#\mathbf{S}(\mathbf{r}, \ell) \leq 2 \binom{N + \lceil |\mathbf{r}|/2 \rceil}{N}.$$

On the other hand, from the generating function  $P_{\mathbf{r}}(t)$  we see that  $c(\mathbf{r}, \ell)$  is bounded by  $\binom{|\mathbf{r}|}{\ell}$  when  $|\mathbf{r}|$  is even and  $\ell = |\mathbf{r}|/2$  and  $2\binom{|\mathbf{r}|}{\ell}$  otherwise. This follows from the fact that when  $\mathbf{k}$  is the multi-index having  $|\mathbf{r}|$  ones and the rest zeros, since  $\ell \leq |\mathbf{r}|$ , we have that  $\#\mathbf{S}(\mathbf{r}, \ell) \leq \#\mathbf{T}(\mathbf{k}, \ell)$  where  $\mathbf{T}(\mathbf{k}, \ell) = \{\mathbf{s} \in \mathbb{N}_0^N : |\mathbf{s}| = \ell, \mathbf{s} \leq \mathbf{k}\}$  and  $\#\mathbf{T}(\mathbf{k}, \ell)$  is given by the coefficient of  $t^\ell$  in  $P_{\mathbf{k}}(t) = (1+t)^{|\mathbf{r}|} = \sum_{\ell=0}^{|\mathbf{r}|} \binom{|\mathbf{r}|}{\ell} t^\ell$  from the binomial theorem. Then

$$\sum_{\ell=\lceil |\mathbf{r}|/2 \rceil}^{|\mathbf{r}|} c(\mathbf{r}, \ell) \leq 2 \sum_{\ell=0}^{|\mathbf{r}|} c(\mathbf{r}, \ell) = 2^{|\mathbf{r}|+1},$$

so that

$$\text{nnz}(\mathbf{G}_{\mathbf{r}}) = \sum_{\ell=\lceil |\mathbf{r}|/2 \rceil}^{|\mathbf{r}|} c(\mathbf{r}, \ell) \binom{N+p-\ell}{p-\ell} \leq 2 \min \left\{ 2^{|\mathbf{r}|}, \binom{N + \lceil |\mathbf{r}|/2 \rceil}{N} \right\} \binom{N+p-\lceil |\mathbf{r}|/2 \rceil}{N},$$

showing (30). Substituting (30) into (26) shows the bound of  $\mathcal{M}(p, r)$  from (31).  $\square$

We note that the bound of  $\mathcal{M}(p, r)$  from (31) is an overestimate due to the particular form of (29), which is different for each  $\mathbf{r} \in \Lambda_r$ . As a consequence, we see that  $\text{nnz}(\mathbf{G}_{\mathbf{r}}) = \mathcal{O}(\min\{2^{|\mathbf{r}|}, M_{\lceil |\mathbf{r}|/2 \rceil}\} M_{p-\lceil |\mathbf{r}|/2 \rceil})$  for  $\mathbf{r} \in \Lambda_r$ . For large  $N$  and small  $r$ ,  $2^r$  is smaller than  $M_{\lceil r/2 \rceil}$ , however, as  $r \rightarrow \infty$  the bound  $M_{\lceil r/2 \rceil}$  is sharper.

For the  $\varepsilon$ -complexity analysis in the next section, we note that

$$\begin{aligned}
\mathcal{M}(p, r) &= \sum_{\mathbf{r} \in \Lambda_r} \text{nnz}(\mathbf{G}_{\mathbf{r}}) \leq \sum_{\mathbf{r} \in \Lambda_r} 2 \min \left\{ 2^{|\mathbf{r}|}, \binom{N + \lceil |\mathbf{r}|/2 \rceil}{N} \right\} \binom{N + p - \lceil |\mathbf{r}|/2 \rceil}{N} \\
&= \sum_{j=0}^r 2 \min \left\{ 2^j, \binom{N + \lceil j/2 \rceil}{N} \right\} \binom{N - 1 + j}{N - 1} \binom{N + p - \lceil j/2 \rceil}{N} \\
&\leq 2 \min \left\{ 2^r, \binom{N + \lceil r/2 \rceil}{N} \right\} \binom{N + p}{N} \sum_{j=0}^r \binom{N - 1 + j}{N - 1} \\
&= 2 \min \left\{ 2^r, \binom{N + \lceil r/2 \rceil}{N} \right\} \binom{N + p}{N} \binom{N + r}{N}
\end{aligned} \tag{33}$$

Figure 4 plots how sharply  $\mathcal{M}(p, r)$  is bounded by (31) and (33). We are also able to show that Theorem 4.1 yields a sharp result in the case  $|\mathbf{r}| = 1$ .

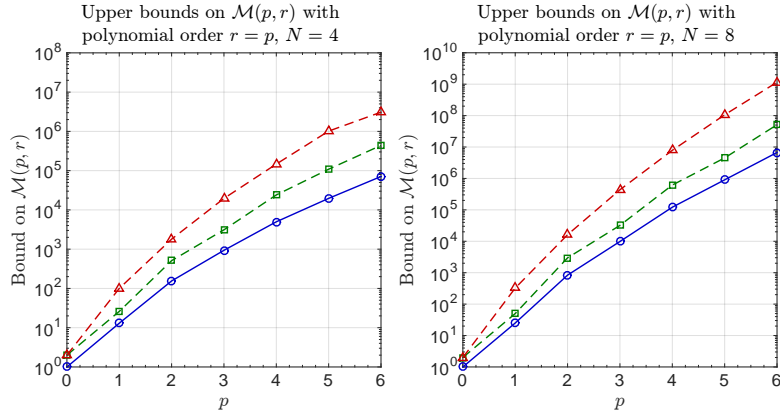


Figure 4: For  $r = p$  with  $p$  ranging from  $0, 1, \dots, 6$  we plot for  $N = 4$  (left) and  $N = 8$  (right) the actual sparsity  $\mathcal{M}(p, r)$  given by (26) of the Galerkin system  $\mathbf{K}_r$  from (22) (blue), the bound on the sparsity from (31) (green), and the bound on the sparsity from (33) (red).

**Corollary 4.2.** *Under the same conditions in Theorem 4.1, when  $\mathbf{r} \in \Lambda_r$  is such that  $|\mathbf{r}| = 1$ , we have*

$$\text{nnz}(\mathbf{G}_{\mathbf{r}}) = \sum_{\ell=\lceil |\mathbf{r}|/2 \rceil}^{|\mathbf{r}|} c(\mathbf{r}, \ell) \binom{N + p - \ell}{p - \ell} = 2 \binom{N + p - 1}{p - 1}. \tag{34}$$

Corollary 4.2 is the result of [14, Lemma 28], and follows from the application of the exact formula for  $\text{nnz}(\mathbf{G}_{\mathbf{r}})$  from (29). Here  $|\mathbf{r}| = \lceil |\mathbf{r}|/2 \rceil = 1$  is odd and  $\mathbf{S}(\mathbf{r}, 1)$  has only one element  $\mathbf{S}(\mathbf{r}, 1) = \{\mathbf{s} \in \mathbb{N}_0^N : |\mathbf{s}| = 1, \mathbf{s} \leq \mathbf{r}\} = \{\mathbf{r}\}$ . Hence  $c(\mathbf{r}, 1) = 2\#\mathbf{S}(\mathbf{r}, 1) = 2$ , and (34) is shown. We are also able to show that the formula for  $\text{nnz}(\mathbf{G}_{\mathbf{r}})$  from (29) yields a result that is sharp in the case  $N = 1$ .

**Corollary 4.3.** *Under the same conditions in Theorem 4.1, when  $N = 1$  and  $\mathbf{r} = r \in \mathbb{N}_0$ , we have*

(a) *in case  $r = 2k$ ,  $k \in \mathbb{N}_0$ ,*

$$\text{nnz}(\mathbf{G}_r) = \begin{cases} (p - r + 1)(r + 1) + k^2, & 0 \leq r \leq p, \\ (p - k + 1)^2, & p + 1 \leq r \leq 2p, \\ 0, & r > 2p. \end{cases} \tag{35}$$

(b) in case  $r = 2k + 1$ ,  $k \in \mathbb{N}_0$ ,

$$\text{nnz}(\mathbf{G}_r) = \begin{cases} (p - r + 1)(r + 1) + k^2 + k, & 0 \leq r \leq p, \\ (p - k + 1)(p - k), & p + 1 \leq r \leq 2p, \\ 0, & r > 2p. \end{cases} \quad (36)$$

Corollary 4.3 is the result of [14, Lemma 25], and its proof using Theorem 4.1 is included in the Appendix. In the remarks that follow, we make a distinction between the cases that  $a(x, \mathbf{y})$  is a polynomial of fixed degree  $\bar{r} < \infty$ , e.g., the coefficients from Examples 2.1 and 2.2, and that  $a(x, \mathbf{y})$  is a transcendental function of the random variables, e.g., the coefficient from Example 2.3.

**Remark 4.4.** (Complexity of matrix-vector products for polynomial coefficients, see e.g., Examples 2.1 and 2.2) From Corollary 4.2 and the work estimate (28) when using (25) as a preconditioner, we see that for coefficients that are affine functions of the random variables, e.g., Example 2.1, the complexity of a single PCG iteration is of the order  $\mathcal{O}(J_h(2M_p + 2NM_{p-1})) = \mathcal{O}(J_h M_p)$ , where  $M_p = \#\Lambda_p^{\text{TD}} = \binom{N+p}{N}$ . On the other hand, when the coefficient  $a(x, \mathbf{y})$  is a polynomial function of the random variables, e.g., Example 2.2, having fixed order  $\bar{r} \in \mathbb{N}$ ,  $\bar{r} < \infty$ , we use Theorem 4.1 to obtain a different estimate. Since  $\{\Psi_{\mathbf{r}}\}_{\mathbf{r} \in \Lambda_{\bar{r}}}$  is a basis for the space  $\mathcal{P}_{\Lambda_{\bar{r}}}(\Gamma)$ , there exists coefficients  $\{a_{\mathbf{r}}(x)\}_{\mathbf{r} \in \Lambda_{\bar{r}}}$  such that  $a(x, \mathbf{y}) = \bar{a}(x, \mathbf{y}) = \sum_{\mathbf{r} \in \Lambda_{\bar{r}}} a_{\mathbf{r}}(x) \Psi_{\mathbf{r}}(\mathbf{y})$ . With this representation, it is clear to see that substituting  $a(x, \mathbf{y})$  into (19) yields  $\mathbf{K}_{\bar{r}}$  from (22), and  $\mathbf{K}_{\bar{r}} = \mathbf{K}$  from (16). However, it is not clear how many of the coefficients  $a_{\mathbf{r}}(x)$  are identically zero. In this case, we can provide an upper bound on the block-sparsity of  $\mathbf{K}_{\bar{r}}$  under the assumption that  $a_{\mathbf{r}}(x) \not\equiv 0 \forall \mathbf{r} \in \Lambda_{\bar{r}}$ . Using the bound of (33), the complexity of a single matrix-vector product of  $\mathbf{K}_{\bar{r}}$  is of the order  $\mathcal{O}(J_h M_p M_{\bar{r}} \min\{2^{\bar{r}}, M_{\lceil \bar{r}/2 \rceil}\})$ . Thus, when  $\bar{r}$  is fixed,  $M_{\bar{r}} \min\{2^{\bar{r}}, M_{\lceil \bar{r}/2 \rceil}\}$  is a constant, and this estimate has the same asymptotic complexity as  $\mathcal{O}(J_h M_p)$ .

**Remark 4.5.** (Complexity of matrix-vector products in the transcendental case, see e.g., Example 2.3) We recall the discussion of [34, Section 3.4]. There, the complexity of matrix-vector products with the SG system was estimated when a full orthogonal expansion is substituted into the SG discretization. This case corresponds to fixing the expansion order  $r = 2p$  following Remark 3.3. Assuming that  $\text{nnz}(\mathbf{G}_r) = \mathcal{O}(M_p)$  or  $\mathcal{O}(M_p^2)$ , it was estimated in [35] that the cost of matrix-vector products involving  $\mathbf{K}_r$  is between  $\mathcal{O}(J_h M_p^2)$  and  $\mathcal{O}(J_h M_p^3)$ . However, the use of Theorem 4.1 allows us to consider the complexity in the case of truncating the expansion, where a sharper estimate can be obtained. Let  $T_r := \prod_{k=\lceil r/2 \rceil+1}^r \frac{N+k}{k} \ll M_{\lceil r/2 \rceil} = \binom{N+\lceil r/2 \rceil}{\lceil r/2 \rceil}$ , which is bounded independent of  $r$ , i.e.,

$$T_r \leq \left( \frac{N + \lceil r/2 \rceil + 1}{\lceil r/2 \rceil + 1} \right)^{\lceil r/2 \rceil + 1} \rightarrow e^N \quad \text{as } r \rightarrow \infty,$$

so that  $M_r = T_r M_{\lceil r/2 \rceil} \leq e^N M_{\lceil r/2 \rceil}$ . From (33), we see that  $\mathcal{M}(p, r)$  is of the order  $\mathcal{O}(M_p M_r M_{\lceil r/2 \rceil})$  as  $p, r \rightarrow \infty$ , since  $\min\{2^r, M_{\lceil r/2 \rceil}\} \rightarrow M_{\lceil r/2 \rceil}$  as  $r \rightarrow \infty$ . When  $r = 2p$ , this implies the complexity of matrix-vector multiplications involving  $\mathbf{K}_r$  is of the order  $\mathcal{O}(J_h M_p M_r M_{\lceil r/2 \rceil}) = \mathcal{O}(J_h M_p T_r M_{\lceil r/2 \rceil}^2) = \mathcal{O}(J_h M_p^3)$ . On the other hand, when  $r = p$ , we see that the complexity of matrix-vector products with  $\mathbf{K}_r$  is order  $\mathcal{O}(J_h M_p M_r M_{\lceil r/2 \rceil}) = \mathcal{O}(J_h T_r^2 M_{\lceil p/2 \rceil}^3) = \mathcal{O}(J_h M_{\lceil p/2 \rceil}^3)$ .

#### 4.2. $\varepsilon$ -complexity analysis of the SGFEM

An estimate of the total complexity to obtain a fully discrete approximation of tolerance  $\varepsilon > 0$  with the SGFEM and PCG solver can be shown in four steps:

1. Estimate the maximum mesh size  $h_{\max}$  and minimum polynomial order  $p_{\min}$  necessary in the finite element and SG discretizations, respectively,
2. If projection of the coefficient is necessary, estimate the minimum projection order  $r_{\min}$ , otherwise set  $r_{\min} = \bar{r}$  where  $\bar{r} < \infty$  is the order of the coefficient,
3. Estimate the minimum number of iterations  $k_{\min}$  needed by the PCG solver,

4. Substitute  $h_{\max}$ ,  $p_{\min}$ ,  $r_{\min}$ , and  $k_{\min}$  into the cost (28) for  $h$ ,  $p$ ,  $r$ , and  $N_{\text{iter}}^{\text{PSG}}$ , respectively.

We proceed to estimate these parameters as follows. Denote by  $u^r$  the corresponding solution of (1) when  $a^r(x, \mathbf{y})$  is substituted in place of  $a(x, \mathbf{y})$ , and let  $\tilde{u}_{h,p}^r$  be the approximation to  $u_{h,p}^r$  found by PCG. Then the total error for the SG approximation satisfies the following bound:

$$\|u - \tilde{u}_{h,p}^r\|_{\mathcal{H}_e^2} \leq \underbrace{\|u - u^r\|_{\mathcal{H}_e^2}}_{\text{SG(I)}} + \underbrace{\|u^r - u_h^r\|_{\mathcal{H}_e^2}}_{\text{SG(II)}} + \underbrace{\|u_h^r - u_{h,p}^r\|_{\mathcal{H}_e^2}}_{\text{SG(III)}} + \underbrace{\|u_{h,p}^r - \tilde{u}_{h,p}^r\|_{\mathcal{H}_e^2}}_{\text{SG(IV)}}. \quad (37)$$

In this setting SG(I) is the approximation error using a truncated expansion of  $a(x, \mathbf{y})$ , SG(II) is the discretization error induced by the finite element method, SG(III) is the SG error coming from the orthogonal expansion (13), and SG(IV) is the solver error resulting from the PCG method. We note that when the projection of the coefficient is exact, as discussed in Remark 3.3, the approximation error SG(I) is no longer present and  $u_{h,p}^r \equiv u_{h,p}$ .

We start with bounding SG(III). Without loss of generality, it is reasonable to assume that since  $u^r$  has a holomorphic dependence on  $\mathbf{z} \in \mathbb{C}^N$  in an open neighborhood of the polyellips  $\mathcal{E}_\gamma$  from Theorem 2.4, then  $u_h^r$  does as well. Then, the following result, whose proof is found in [32], and immediately follows from classical spectral convergence results [9, 31], describes the convergence rate of the fully discrete solutions obtained by the SG method using a total degree approximation in  $\mathcal{P}_{\Lambda_p^{\text{TD}}}(\Gamma)$ :

**Proposition 4.6** (Convergence rate for the SG method). *If Theorem 2.4 holds for the solution  $u_h^r$  to (9) with coefficient  $a^r(x, \mathbf{y})$ , and  $u_{h,p}^r$  is the solution to (14) with  $\Lambda_p$  the order  $p$  total degree index set, then*

$$\|u_h^r - u_{h,p}^r\|_{\mathcal{H}_e^\infty} \leq C_1 \exp(-C_2 p) \quad \forall p \in \mathbb{N},$$

for some constants  $C_1, C_2 > 0$  independent of  $p$ .

To investigate the error in SG(I), we note that since  $a(x, \mathbf{y})$  satisfies assumption (A2), the projection error in  $\mathcal{P}_{\Lambda_p^{\text{TD}}}(\Gamma)$  can be similarly estimated as

$$\|a - a^r\|_{L_e^2(\Gamma; L^\infty(D))} \leq C_3 \exp(-C_4 r) \quad \forall r \in \mathbb{N}, \quad (38)$$

for some constants  $C_3, C_4 > 0$  independent of  $r$ . Hence,  $\forall r \in \mathbb{N}$ ,

$$\|u - u^r\|_{\mathcal{H}_e^2} \leq \frac{\|f\|_{H^{-1}}}{a_{\min}^2} \|a - a^r\|_{L_e^2(\Gamma; L^\infty(D))} \leq \frac{\|f\|_{H^{-1}}}{a_{\min}^2} C_3 \exp(-C_4 r) \quad (39)$$

providing a bound for SG(I). For a bound of SG(II), we present the following convergence result regarding solutions to the parameterized finite element problem, whose proof can be found in a number of standard texts on the theory of finite element methods, e.g., [2, 20]:

**Lemma 4.7.** *Let  $\mathcal{T}_h$  be a uniform finite element mesh over  $D$  with  $J_h = \mathcal{O}(h^{-d})$  degrees of freedom and  $h > 0$ . For the elliptic PDE (1) and  $\mathbf{y} \in \Gamma$ , when  $u^r(\mathbf{y}) \in H_0^1(D) \cap H^{s+1}(D)$ , the error from the finite element approximation is bounded by*

$$\|u^r(\mathbf{y}) - u_h^r(\mathbf{y})\|_{H_0^1(D)} \leq C_{\text{FEM}} h^s,$$

where the constant  $C_{\text{FEM}} > 0$  is independent of  $h$  and  $\mathbf{y}$ .

For the treatment of SG(IV), we begin by defining  $\mathcal{B}^r(\mathbf{y})$  to be the corresponding bilinear operator in (6) with  $a(x, \mathbf{y})$  replaced with  $a^r(x, \mathbf{y})$ . Since both  $\mathcal{B}(\mathbf{y})$  and  $\mathcal{B}^r(\mathbf{y})$  are symmetric, uniformly coercive and continuous bilinear operators on  $H_0^1(D)$ , there exist  $\alpha, \beta > 0$  independent of  $\mathbf{y}$  such that for every  $u, v \in H_0^1(D)$

$$\begin{aligned} |\mathcal{B}^r[u, v](\mathbf{y})| &= \left| \int_D a^r(x, \mathbf{y}) \nabla u \cdot \nabla v dx \right| \leq \alpha \|u\|_{H_0^1(D)} \|v\|_{H_0^1(D)}, \quad \text{and} \\ \beta \|u\|_{H_0^1(D)}^2 &\leq \int_D a^r(x, \mathbf{y}) |\nabla u|^2 dx = \|u\|_{\mathcal{B}^r(\mathbf{y})}^2, \end{aligned}$$



and similarly for  $\mathcal{B}_r(\mathbf{y})$  with the same  $\alpha, \beta$ , e.g., taking  $\alpha$  to be the maximum and  $\beta$  to be the minimum in each case. Recall  $\mathbf{u}_{h,p}^r = [u_{1,p}^r, \dots, u_{J_h,p}^r]^T$ , the vector of nodal values of the finite element solution corresponding to the  $p$ -th stochastic mode of  $u_{h,p}^r$ , and  $\mathbf{u}_{h,p}^r = [\mathbf{u}_{h,p}^r]_{\mathbf{p} \in \Lambda_p}^T$ . Then we have the following estimates expressing

$$\text{Continuity:} \quad \|\mathbf{u}_{h,p}^r\|_{\mathbf{K}_r} = \|u_{h,p}^r\|_{\mathbb{E}[\mathcal{B}^r(\mathbf{y})]} \leq \sqrt{\alpha} \|u_{h,p}^r\|_{\mathcal{H}_\varepsilon^2}, \quad \text{and} \quad (40)$$

$$\text{Ellipticity:} \quad \sqrt{\beta} \|u_{h,p}^r\|_{\mathcal{H}_\varepsilon^2} \leq \|u_{h,p}^r\|_{\mathbb{E}[\mathcal{B}^r(\mathbf{y})]} = \|\mathbf{u}_{h,p}^r\|_{\mathbf{K}_r}, \quad (41)$$

where  $\|\mathbf{u}\|_{\mathbf{K}_r}^2 = (\mathbf{u})^T \mathbf{K}_r \mathbf{u}$  is the  $\mathbf{K}_r$  matrix norm, and  $\|u\|_{\mathbb{E}[\mathcal{B}^r(\mathbf{y})]}$  is the expectation of the energy norm (7). Given Proposition 4.6, Lemma 4.7, and the estimates from (39), (40), and (41), we can now provide the minimal projection orders  $p, r \in \mathbb{N}$  for the SG approximation (13) and the coefficient (18), respectively, the maximum mesh size  $h$  for finite element method, and the minimum number of PCG iterations  $k$  necessary to ensure that the error in the SGFEM solution  $\tilde{u}_{h,p}^r$  is less than the tolerance  $\varepsilon > 0$ .

**Lemma 4.8.** *Let  $u \in L_\varepsilon^2(\Gamma; H_0^1(D) \cap H^{s+1}(D))$  be the solution to (1),  $u_{h,p}^r$  be the solution to (14) with the coefficient  $a^r(x, \mathbf{y})$ , and  $\tilde{u}_{h,p}^r$  be the approximation of  $u_{h,p}^r$  found by PCG with a zero initial guess. Then, for  $\varepsilon > 0$ , to ensure that  $\|u - \tilde{u}_{h,p}^r\|_{\mathcal{H}_\varepsilon^2} \leq \varepsilon$  we must choose  $h \leq h_{\max}$ ,  $r \geq r_{\min}$ ,  $p \geq p_{\min}$ , and  $k \geq k_{\min}$ , where:*

$$\begin{aligned} h_{\max} &= \left( \frac{\varepsilon}{4C_{\text{FEM}}} \right)^{\frac{1}{s}}, & r_{\min} &= \log \left[ \left( \frac{4C_5}{\varepsilon} \right)^{\frac{1}{C_4}} \right], \\ p_{\min} &= \log \left[ \left( \frac{4C_1}{\varepsilon} \right)^{\frac{1}{C_2}} \right], & k_{\min} &= \frac{\log \left( \frac{4C_6}{\varepsilon} \right)}{\log \left( \frac{\sqrt{\tilde{\kappa}_r} + 1}{\sqrt{\tilde{\kappa}_r} - 1} \right)}, \end{aligned}$$

with  $C_{\text{FEM}} > 0$  the constant from Lemma 4.7,  $C_1, C_2, C_3, C_4 > 0$  the constants from Proposition 4.6 and (38), and, for  $\alpha, \beta > 0$  from (40) and (41)

$$C_5 = C_3 \frac{\|f\|_{H^{-1}}}{a_{\min}^2}, \quad C_6 = 2\sqrt{\frac{\alpha}{\beta}} \|u_{h,p}^r\|_{\mathcal{H}_\varepsilon^2},$$

with  $\tilde{\kappa}_r$  is the condition number of  $\tilde{\mathbf{P}}^{-1} \mathbf{K}_r$  with  $\mathbf{P}$  the mean-based preconditioner from (25).

**PROOF.** Without loss of generality, we seek to bound the quantities SG(I)-SG(IV) from (37) each by  $\varepsilon/4$ . For the error SG(I) we recall estimate (39) and solve for  $r$ . From Lemma 4.7, when  $u \in L_\varepsilon^2(\Gamma; H_0^1(D) \cap H^{s+1}(D))$  we have that  $\|u^r - u_h^r\|_{\mathcal{H}_\varepsilon^2} \leq C_{\text{FEM}} h^s \forall h > 0$ , and from Proposition 4.6 we have that  $\|u_h^r - u_{h,p}^r\|_{\mathcal{H}_\varepsilon^\infty} \leq C_1 \exp(-C_2 p) \forall p \in \mathbb{N}$ , so that solving for  $h$  and  $p$  gives the desired maximum mesh size  $h_{\max}$  and minimum polynomial order  $p_{\min}$  to bound SG(II) and SG(III) by  $\varepsilon/4$ . Let  $\mathbf{u}_{h,p}^r$  and  $\mathbf{u}_{h,p}^{r,(k)}$  be the coefficients of the exact SG solution  $u_{h,p}^r$  and the approximate SG solution  $\tilde{u}_{h,p}^r$  after  $k$  PCG iterations, respectively. Then from (24) and (41) we see that

$$\|u_{h,p}^r - \tilde{u}_{h,p}^r\|_{\mathcal{H}_\varepsilon^2} \leq \frac{1}{\sqrt{\beta}} \|\mathbf{u}_{h,p}^r - \mathbf{u}_{h,p}^{r,(k)}\|_{\mathbf{K}_r} \leq \frac{2}{\sqrt{\beta}} \left( \frac{\sqrt{\tilde{\kappa}_r} - 1}{\sqrt{\tilde{\kappa}_r} + 1} \right)^k \|\mathbf{u}_{h,p}^r - \mathbf{u}_{h,p}^{r,(0)}\|_{\mathbf{K}_r},$$

where  $\mathbf{u}_{h,p}^{r,(0)}$  is the initial guess used in CG and  $\tilde{\kappa}_r = \text{cond}(\tilde{\mathbf{P}}^{-1} \mathbf{K}_r)$  with mean based preconditioner  $\mathbf{P}$  from (25). If we use the zero vector as the initial iteration in PCG, we have from (40)

$$\|u_{h,p}^r - \tilde{u}_{h,p}^r\|_{\mathcal{H}_\varepsilon^2} \leq \frac{2}{\sqrt{\beta}} \left( \frac{\sqrt{\tilde{\kappa}_r} - 1}{\sqrt{\tilde{\kappa}_r} + 1} \right)^k \|\mathbf{u}_{h,p}^r\|_{\mathbf{K}_r} \leq 2\sqrt{\frac{\alpha}{\beta}} \left( \frac{\sqrt{\tilde{\kappa}_r} - 1}{\sqrt{\tilde{\kappa}_r} + 1} \right)^k \|u_{h,p}^r\|_{\mathcal{H}_\varepsilon^2}. \quad (42)$$

Solving for  $k$  gives the minimum number of iterations  $k_{\min}$  required to ensure SG(IV) is bounded by  $\varepsilon/4$ .  $\square$

Given the necessary parameters from Lemma 4.8 to achieve  $\|u - \tilde{u}_{h,p}^r\|_{\mathcal{H}_e^2} \leq \varepsilon$ , and the estimates on the computational complexity of one iteration in the PCG method from §4.1, we provide a bound on the minimal number of FLOPs required by the SGFEM when approximating (1). We split these results into the cases that the stochastic coefficient  $a(x, \mathbf{y})$  from (1) is:

- (i) an affine function of the random parameters, e.g.,  $a(x, \mathbf{y}) \in \mathcal{P}_{\Lambda_{\bar{r}}}(\Gamma)$  with  $\bar{r} = 1$ , as in Example 2.1,
- (ii) a non-affine polynomial of the random parameters, e.g.,  $a(x, \mathbf{y}) \in \mathcal{P}_{\Lambda_{\bar{r}}}(\Gamma)$  for some  $1 < \bar{r} < \infty$ , as in Example 2.2,
- (iii) a non-affine, transcendental function of the random parameters, e.g.,  $a(x, \mathbf{y}) \notin \mathcal{P}_{\Lambda_{\bar{r}}}(\Gamma)$  for any  $r \in \mathbb{N}$ , as in Example 2.3, so that  $r$  must be chosen to satisfy  $r \geq r_{\min}$  from Lemma 4.8.

The results are summarized in Theorems 4.9 and 4.10 next.

**Theorem 4.9.** *Let  $u \in L^2_\rho(\Gamma; H^1_0(D) \cap H^{s+1}(D))$  be the solution to (1), and  $\bar{r}$  be the smallest natural number such that  $a(x, \mathbf{y}) \in \mathcal{P}_{\Lambda_{\bar{r}}}(\Gamma)$ . When  $\bar{r} = 1$ , the minimum work (28) of solving (22) with PCG to a tolerance  $\varepsilon > 0$  can be bounded by*

$$W^{\text{pSG}} \leq C_7 \left( \frac{3C_{\text{FEM}}}{\varepsilon} \right)^{\frac{d}{s}} 2e^N (1+N) \left( 1 + \log \left[ \left( \frac{3C_1}{\varepsilon} \right)^{\frac{1}{C_2 N}} \right] \right)^N \left( \frac{\log \left( \frac{3C_6}{\varepsilon} \right)}{\log \left( \frac{\sqrt{\kappa}+1}{\sqrt{\kappa}-1} \right)} \right), \quad (43)$$

and when  $\bar{r} > 1$ , the minimum work (28) of solving (22) with PCG to a tolerance  $\varepsilon > 0$  can be bounded by

$$\begin{aligned} W^{\text{pSG}} &\leq C_7 \left( \frac{3C_{\text{FEM}}}{\varepsilon} \right)^{\frac{d}{s}} \left( \frac{\log \left( \frac{3C_6}{\varepsilon} \right)}{\log \left( \frac{\sqrt{\kappa}+1}{\sqrt{\kappa}-1} \right)} \right) e^N \left[ \left( 1 + \log \left[ \left( \frac{3C_1}{\varepsilon} \right)^{\frac{1}{C_2 N}} \right] \right)^N \right. \\ &\quad \left. + 2 \sum_{j=0}^{\bar{r}} \binom{N-1+j}{N-1} \min \left\{ 2^j, \binom{N+\lceil j/2 \rceil}{N} \right\} \left( 1 - \frac{\lceil j/2 \rceil}{N} + \log \left[ \left( \frac{3C_1}{\varepsilon} \right)^{\frac{1}{C_2 N}} \right] \right)^N \right], \end{aligned} \quad (44)$$

with  $C_{\text{FEM}}, C_1, C_2, C_6$  as in Lemma 4.8,  $C_7 > 0$  independent of  $\varepsilon$ , and  $\kappa$  the condition number of the preconditioned system  $\tilde{\mathbf{P}}^{-1} \mathbf{K}_{\bar{r}} = \tilde{\mathbf{P}}^{-1} \mathbf{K}$ , using the mean-based preconditioner from (25).

PROOF. When  $a(x, \mathbf{y}) \in \mathcal{P}_{\Lambda_{\bar{r}}}(\Gamma)$  we do not need to consider SG(I) from (37), and bound SG(II), SG(III), and SG(IV) by  $\varepsilon/3$ . Hence, to minimize the error of the SG discretization, we choose  $p \geq p_{\min} = \log[(3C_1/\varepsilon)^{1/C_2}]$  which differs from the  $p_{\min}$  stated in Lemma 4.8. For a uniform triangulation  $\mathcal{T}_h$ ,  $J_h = \mathcal{O}(h^{-d})$  so that

$$J_{h_{\max}} = C_7 \left[ \left( \frac{\varepsilon}{3C_{\text{FEM}}} \right)^{\frac{1}{s}} \right]^{-d} = C_7 \left( \frac{3C_{\text{FEM}}}{\varepsilon} \right)^{\frac{d}{s}} \quad (45)$$

for some constant  $C_7 > 0$  depending on the connectivity of the finite element mesh, but independent of  $\varepsilon$ . In the case that  $\bar{r} = 1$ , we substitute  $p_{\min}$  into (29) for the matrices  $\mathbf{G}_{\mathbf{r}}$  having  $0 \leq |\mathbf{r}| \leq 1$ , and apply Stirling's approximation to obtain

$$M_{p_{\min}} + \mathcal{M}(p_{\min}, 1) = 2 \binom{N+p_{\min}}{N} + 2N \binom{N+p_{\min}-1}{N} \leq 2e^N (1+N) \left( 1 + \log \left[ \left( \frac{3C_1}{\varepsilon} \right)^{\frac{1}{C_2 N}} \right] \right)^N,$$

Similarly, when  $\bar{r} > 1$  we use the bound from (31) and Stirling's approximation to obtain

$$\begin{aligned} M_{p_{\min}} + \mathcal{M}(p_{\min}, \bar{r}) &\leq e^N \left[ \left( 1 + \log \left[ \left( \frac{3C_1}{\varepsilon} \right)^{\frac{1}{C_2 N}} \right] \right)^N \right. \\ &\quad \left. + 2 \sum_{j=0}^{\bar{r}} \min \left\{ 2^j, \binom{N+\lceil j/2 \rceil}{N} \right\} \binom{N-1+j}{N-1} \left( 1 - \frac{\lceil j/2 \rceil}{N} + \log \left[ \left( \frac{3C_1}{\varepsilon} \right)^{\frac{1}{C_2 N}} \right] \right)^N \right]. \end{aligned}$$

Substituting  $J_{h_{\max}}$  for  $J_h$ ,  $k_{\min}$  for  $N_{\text{iter}}^{\text{SG}}$  from Lemma 4.8, and the bounds for  $M_{p_{\min}} + N_K^{(p_{\min}, \bar{r})}$  into the work estimate (28), in the cases  $\bar{r} = 1$  and  $\bar{r} > 1$  above, we obtain the desired results.  $\square$

**Theorem 4.10.** *Let  $u \in L^2_\rho(\Gamma; H^1_0(D) \cap H^{s+1}(D))$  be the solution to (1), and suppose that  $a(x, \mathbf{y}) \notin \mathcal{P}_{\Lambda_r}(\Gamma)$  for any  $r \in \mathbb{N}$ . In this case  $r$  must be chosen to satisfy  $r \geq r_{\min}$  from Lemma 4.8. Then the minimum work (28) of solving (22) with PCG to a tolerance  $\varepsilon > 0$  can be bounded by*

$$W^{\text{PSG}} \leq C_7 \left( \frac{4C_{\text{FEM}}}{\varepsilon} \right)^{\frac{d}{s}} e^N \left( 1 + \log \left[ \left( \frac{4C_1}{\varepsilon} \right)^{\frac{1}{C_2 N}} \right] \right)^N \left( 1 + 2e^{2N} \left( 1 + \log \left[ \left( \frac{4C_5}{\varepsilon} \right)^{\frac{1}{C_4 N}} \right] \right)^N \left( 1 + \frac{1}{N} + \log \left[ \left( \frac{4C_5}{\varepsilon} \right)^{\frac{1}{2C_4 N}} \right] \right)^N \right) \left( \frac{\log \left( \frac{4C_6}{\varepsilon} \right)}{\log \left( \frac{\sqrt{\tilde{\kappa}_r} + 1}{\sqrt{\tilde{\kappa}_r} - 1} \right)} \right), \quad (46)$$

with  $C_{\text{FEM}}, C_1, C_2, C_4, C_5, C_6$  as in Lemma 4.8,  $C_7 > 0$  independent of  $\varepsilon$ , and  $\tilde{\kappa}_r$  the condition number of the preconditioned system  $\tilde{\mathbf{P}}^{-1} \mathbf{K}_r$ , using the mean-based preconditioner from (25).

PROOF. In this setting  $r$  must be chosen to satisfy  $r \geq r_{\min} = \log[(4C_5/\varepsilon)^{1/C_4}]$  from Lemma 4.8 and, therefore, we must bound the sum from (31) which now depends on  $r_{\min}$ , and hence on  $\varepsilon$ . Thus, we use the bound (33) for  $\mathcal{M}(p, r)$ , noting that as  $\varepsilon \rightarrow 0$ ,  $r_{\min} \rightarrow \infty$  so that

$$\min \left\{ 2^{r_{\min}}, \binom{N + \lceil r_{\min}/2 \rceil}{N} \right\} = \binom{N + \lceil r_{\min}/2 \rceil}{N}.$$

Substituting  $p_{\min}$  and  $r_{\min}$  from Lemma 4.8 into (33) and applying Stirling's approximation, we obtain

$$M_{p_{\min}} + \mathcal{M}(p_{\min}, r_{\min}) \leq e^N \left( 1 + \log \left[ \left( \frac{4C_1}{\varepsilon} \right)^{\frac{1}{C_2 N}} \right] \right)^N \left( 1 + 2e^{2N} \left( 1 + \log \left[ \left( \frac{4C_5}{\varepsilon} \right)^{\frac{1}{C_4 N}} \right] \right)^N \left( 1 + \frac{1}{N} + \log \left[ \left( \frac{4C_5}{\varepsilon} \right)^{\frac{1}{2C_4 N}} \right] \right)^N \right).$$

As in the proof of Theorem 4.9, we substitute  $J_{h_{\max}}$  for  $J_h$  from (45),  $k_{\min}$  for  $N_{\text{iter}}^{\text{SG}}$ , and the bound for  $M_{p_{\min}} + \mathcal{M}(p_{\min}, r_{\min})$  with  $p_{\min}$  and  $k_{\min}$  from Lemma 4.8 into the cost (28) to complete the proof.  $\square$

Given Theorems 4.9 and 4.10 we see that the work of obtaining the fully discrete approximation using the SGFEM, with PCG as a solver, is asymptotically given by:

$$\underbrace{\mathcal{O} \left( \frac{1}{\varepsilon} \right)^{\frac{d}{s}}}_{(\text{SG.1})} \underbrace{\left[ \log \left( \frac{1}{\varepsilon} \right) \right]^{g(N)}}_{(\text{SG.2})} \underbrace{\left( \frac{\log \left( \frac{1}{\varepsilon} \right)}{\log \left( \frac{\sqrt{\tilde{\kappa}_r} + 1}{\sqrt{\tilde{\kappa}_r} - 1} \right)} \right)}_{(\text{SG.3})}, \quad (47)$$

where  $g(N) = N$  and  $\tilde{\kappa}_r = \tilde{\kappa}$  if  $a(x, \mathbf{y})$  is an affine or non-affine, polynomial function of the random parameters of fixed order  $\bar{r} < \infty$ , e.g., Examples 2.1 and 2.2, and  $g(N) = 3N$  when  $a(x, \mathbf{y})$  is a non-affine, transcendental function of the random parameters, e.g., Example 2.3, requiring a total degree orthogonal expansion of order  $r \geq r_{\min}$  depending on  $\varepsilon$ . Here, (SG.1), (SG.2), and (SG.3) correspond to the work required by the finite element, SG, and PCG methods, respectively. In particular, (SG.2) corresponds to the estimates for the sparsity of the Galerkin system  $\mathbf{K}_r$  from (22), and represents the number of coupled finite element systems that must be solved simultaneously by the PCG method. However, due to the bound (31), the asymptotic complexity in the cases that  $a(x, \mathbf{y})$  is affine or polynomial in  $\mathbf{y}$  are the same. This does not imply that there is no need to consider the work estimates in these cases separately. Indeed, if  $a(x, \mathbf{y})$  is a polynomial having the representation  $\sum_{\mathbf{r} \in \Lambda_{\bar{r}}} a_{\mathbf{r}}(x) \Psi_{\mathbf{r}}(\mathbf{y})$  where  $a_{\mathbf{r}}(x) \neq 0$  for all  $\mathbf{r} \in \Lambda_{\bar{r}}$ , then the complexity of matrix-vector multiplications with  $\mathbf{K}_r$  is of the order  $\mathcal{O}(J_h M_p M_{\bar{r}} \min\{2^{\bar{r}}, M_{\lceil \bar{r}/2 \rceil}\})$ . Here, the constant  $M_{\bar{r}} \min\{2^{\bar{r}}, M_{\lceil \bar{r}/2 \rceil}\}$  grows rapidly with  $\bar{r}$ , suggesting that higher order polynomial functions of  $\mathbf{y}$  require additional cost.

#### 4.3. Conditioning of the generalized SG system

In this section, we discuss issues related to the conditioning of the linear system that results from the SGFEM discretization. We first recall [14, Theorem 10]: the eigenvalues of the matrices  $\{\mathbf{G}_{\mathbf{r}}\}_{\mathbf{r} \in \Lambda_r}$  from (21) lie in the interval  $[\xi_{\mathbf{r}}, \Xi_{\mathbf{r}}]$ , where

$$\xi_{\mathbf{r}} := \min\{\Psi_{\mathbf{r}}(\mathbf{y}) : \mathbf{y} \in \mathcal{G}^{\mathbf{m}(\mathbf{l})}\}, \quad \Xi_{\mathbf{r}} := \max\{\Psi_{\mathbf{r}}(\mathbf{y}) : \mathbf{y} \in \mathcal{G}^{\mathbf{m}(\mathbf{l})}\}, \quad (48)$$

$\mathcal{G}^{\mathbf{m}(\mathbf{l})}$  is a tensor product grid of Gauss-Legendre quadrature points having  $\mathbf{m}(\mathbf{l}) = (m(l_1), \dots, m(l_N))$  points in each direction, and  $\mathbf{l}$  is such that  $m(l_n) := p + \lceil \frac{k_n+1}{2} \rceil$ ,  $n = 1, \dots, N$ . Since  $a^r(x, \mathbf{y})$  satisfies (A1), the analysis of [29, Theorem 3.8] shows that the eigenvalues for the preconditioned system  $\mathbf{P}^{-1}\mathbf{K}_{\mathbf{r}}$  lie in the interval  $[1 - \underline{\tau}_r, 1 + \bar{\tau}_r]$  where

$$\underline{\tau}_r = \frac{1}{a_{\min}} \sum_{\substack{\mathbf{r} \in \Lambda_r \\ |\mathbf{r}| \neq 0}} \xi_{\mathbf{r}} \|a_{\mathbf{r}}(x)\|_{L^\infty(D)}, \quad \bar{\tau}_r = \frac{1}{a_{\min}} \sum_{\substack{\mathbf{r} \in \Lambda_r \\ |\mathbf{r}| \neq 0}} \Xi_{\mathbf{r}} \|a_{\mathbf{r}}(x)\|_{L^\infty(D)}. \quad (49)$$

As a result of (49), we see that in the case that the projection order  $r$  of the coefficient  $a^r(x, \mathbf{y})$  depends on  $\varepsilon$ , the condition number of the preconditioned system  $\mathbf{P}^{-1}\mathbf{K}_{\mathbf{r}}$  does as well through the number of terms in  $\underline{\tau}_r$  and  $\bar{\tau}_r$ . This should come as no surprise since even in the case of the Karhunen-Loève expansion, the condition number of  $\mathbf{P}^{-1}\mathbf{K}$  depends on the number of terms in the truncated Karhunen-Loève expansion which is chosen a-priori to minimize the error.

## 5. Comparison with the SCFEM

In this section we compare our explicit cost bounds for the SGFEM with the complexity estimates for the SCFEM developed in [16], when solving (1). The basic idea behind the SCFEM is to construct a fully discrete approximation in a subspace of  $V_h(D) \otimes L^2_\rho(\Gamma)$  by collocating semi-discrete solutions  $u_h$  from (9) on a deterministic set of points to obtain solutions  $\{u_h(\cdot, \mathbf{y}_k)\}_{k=1}^{M_L} \in V_h(D)$ .

### 5.1. A generalized SCFEM using Lagrange interpolation

To construct the stochastic collocation (SC) approximation, we consider a class of multi-index sets defined in terms of increasing functions  $\mathbf{m} : \mathbb{N}_+^N \rightarrow \mathbb{N}_+^N$  and  $g : \mathbb{N}_+^N \rightarrow \mathbb{N}_+$ . By  $\mathbf{m}$  we specify the multivariate function  $\mathbf{m}(\mathbf{l}) := (m_1(l_1), \dots, m_N(l_N))$  where each  $m_n : \mathbb{N}_+ \rightarrow \mathbb{N}_+$  is an increasing function, possibly different for each  $n = 1, \dots, N$ . Here the  $m_n$  are referred to as growth functions, specifying how the number of points grows in the direction  $n$ . Associated with  $m_n$  we define the *left-inverse*  $m_n^\dagger : \mathbb{N}_+ \rightarrow \mathbb{N}_+$  by  $m_n^\dagger(q) = \min\{k \in \mathbb{N}_+ : m_n(k) \geq q\}$ , and let  $\mathbf{m}^\dagger(\mathbf{q}) = (m_1^\dagger(q_1), \dots, m_N^\dagger(q_N))$ . In this case, we note that  $m_n^\dagger(m_n(k)) = k$  and  $m_n(m_n^\dagger(k)) \geq k$  for each  $k \in \mathbb{N}_+$  and  $n = 1, \dots, N$ . Given  $\mathbf{m}$  and  $g$  we can define the multi-index set

$$\Lambda_L^{\mathbf{m}, g} = \{\mathbf{q} \in \mathbb{N}_+^N : g(\mathbf{m}^\dagger(\mathbf{q} + \mathbf{1})) \leq L\}, \quad (50)$$

to be used in constructing polynomial approximations. In particular, setting  $m_n(j) = j$  for all  $j \in \mathbb{N}_+$  and  $n = 1, \dots, N$ , and defining

$$g_{\text{TP}}(\mathbf{p}) = \max_{1 \leq n \leq N} p_n, \quad g_{\text{TD}}(\mathbf{p}) = \sum_{n=1}^N (p_n - 1), \quad g_{\text{SM}}(\mathbf{p}) = \sum_{n=1}^N f(p_n), \quad (51)$$

where  $f(p)$  is given in (12), and using the definition of  $\Lambda_L^{\mathbf{m}, g}$  from (50), we obtain the TP, TD, and SM index sets  $\Lambda_L^{\text{TP}}$ ,  $\Lambda_L^{\text{TD}}$ , and  $\Lambda_L^{\text{SM}}$ , respectively, given in (12).

We introduce a sequence of one-dimensional Lagrange interpolation operators  $\mathcal{U}^{m_n(l_n)} : C^0(\Gamma_n) \rightarrow \mathcal{P}_{m_n(l_n)-1}(\Gamma_n)$ . Then for  $v \in C^0(\Gamma)$  the generalized multi-dimensional approximation operator  $\mathcal{I}_L^{\mathbf{m},g} : C^0(\Gamma) \rightarrow \mathcal{P}_{\Lambda_L^{\mathbf{m},g}}(\Gamma)$  is given by

$$\mathcal{I}_L^{\mathbf{m},g}[v](\mathbf{y}) = \sum_{g(\mathbf{l}) \leq L} \sum_{\mathbf{i} \in \{0,1\}^N} (-1)^{|\mathbf{i}|} \left( \bigotimes_{n=1}^N \mathcal{U}_n^{m_n(l_n-i_n)} \right) [v](\mathbf{y}). \quad (52)$$

Construction of the approximation  $\mathcal{I}_L^{\mathbf{m},g}[v](\mathbf{y})$  requires the independent evaluation of samples  $v(\mathbf{y})$  on a deterministic set of distinct collocation points  $\mathcal{G}_L^{\mathbf{m},g}$  having cardinality  $M_L = \#\mathcal{G}_L^{\mathbf{m},g}$ . Applying  $\mathcal{I}_L^{\mathbf{m},g}[\cdot]$  from (52) to the semi-discrete solution  $u_h(x, \mathbf{y})$  of problem (9), we obtain the fully discrete SC approximation

$$u_{h,L}(x, \mathbf{y}) = \mathcal{I}_L^{\mathbf{m},g}[u_h](x, \mathbf{y}). \quad (53)$$

*One-dimensional abscissas.* In this effort, we use three examples for constructing the fully discrete approximation. The first is that of a fully-nested rule constructed on the Clenshaw-Curtis choice of abscissas [5] with function  $g_{\text{TD}}(\mathbf{p})$  from (51) and an isotropic growth rule  $\mathbf{m} = (m, \dots, m)$  with  $m$  given by

$$m(1) = 1, \quad m(l_n) = 2^{l_n-1} + 1 \quad \text{for } l_n > 1, \quad (54)$$

This is the classical Smolyak sparse-tensorization construction [30], and here the choice of  $\mathbf{m}$  corresponds to a doubling growth rule that leads to a nested sequence of multi-dimensional grids, e.g.,  $\mathcal{G}_L^{\mathbf{m},g} \subset \mathcal{G}_{L+1}^{\mathbf{m},g}$ . On the other hand, we can construct a sparse-Smolyak approximation on the Gauss-Legendre abscissas corresponding to the zeros of the Legendre polynomials  $\{\Psi_{\mathbf{p}}\}$ , as defined in §3. When the points are grown isotropically according to the linear growth rule with  $\mathbf{m} = (m, \dots, m)$  and  $m$  defined as

$$m(l_n) = l_n \quad \text{for } l_n \in \mathbb{N}, \quad (55)$$

and  $g_{\text{TD}}(\mathbf{p})$  from (51), we obtain a grid that is not nested. Another construction that yields a sequence of nested grids is that based on the Leja points, defined as the sequence of points satisfying  $y_{k+1} := \arg\max_{y \in \Gamma_n} \prod_{j=1}^k |y - y_j|$  (see [10]). Here we take the Leja sequence of points with  $g_{\text{TD}}$  from (51) and the isotropic linear growth function  $\mathbf{m}$  from (55).

### 5.2. Cost of solving the SCFEM systems

To construct the fully discrete approximation with the SCFEM, we must solve  $M_L$  distinct decoupled finite element systems, each dependent on a realization of the parameters  $\mathbf{y}_k \in \mathcal{G}_L^{\mathbf{m},g}$  for  $k = 1, \dots, M_L$ . Similar to the SGFEM, we can apply the PCG method to the solution of each system. Let  $N_{\text{iter}}^{(k)}$  be the number of iterations required by the CG method to solve the finite element system corresponding to  $\mathbf{y}_k$  and  $N_{\text{iter}}^{\text{p}(k)}$  be the corresponding number of iterations when a preconditioner is used. We are interested in choosing a suitable preconditioning strategy to decrease the total number of iterations  $N_{\text{iter}}^{\text{pSC}} = \sum_{k=1}^{M_L} N_{\text{iter}}^{\text{p}(k)}$  required to obtain the fully discrete approximation  $u_{h,L}$ . We present a preconditioning strategy of choosing

$$\mathbf{P}_0 := \mathbf{A}(\mathbf{y}_1), \quad (56)$$

with  $\mathbf{A}(\mathbf{y})$  from (10) the finite element stiffness matrix corresponding to the sample point  $\mathbf{y}_1 \in \mathcal{G}_L^{\mathbf{m},g}$ , as the preconditioner for all of the individual finite element solutions. We refer to this choice of preconditioner as the level-zero preconditioner since it corresponds to the SC approximation at level  $L = 0$ .

Since we apply CG to the solution of each individual finite element system, the work in floating point operations (FLOPs) required to obtain a fully discrete approximation with the SCFEM without a preconditioner is given by

$$W^{\text{SC}} \approx \mathcal{O}(J_h) * \sum_{k=1}^{M_L} N_{\text{iter}}^{(k)}. \quad (57)$$

On the other hand, the “level-zero” preconditioner induces an additional matrix-vector product requiring  $\mathcal{O}(J_h)$  FLOPs per iteration when a sparse factorization of  $\mathbf{P}_0$  is used. Hence the work of solving (5) with PCG is given by

$$W^{\text{pSC}} \approx 2 * \mathcal{O}(J_h) * \sum_{k=1}^{M_L} N_{\text{iter}}^{\text{p}(k)}. \quad (58)$$

Here, the reduction in work due to preconditioning will be seen in the number of iterations saved in each individual count  $N_{\text{iter}}^{\text{p}(k)}$  contributing to the sum.

### 5.3. Comparing the explicit cost bounds of the SGFEM and SCFEM

Given a particular “sparse” index set  $\Lambda_p$ , we can find increasing functions  $\mathbf{m} : \mathbb{N}_+^N \rightarrow \mathbb{N}_+^N$  and  $g : \mathbb{N}_+^N \rightarrow \mathbb{N}_+$ , and  $L \in \mathbb{N}$  such that  $\Lambda_p = \Lambda_L^{\mathbf{m},g}$  from (50). In this setting, we can either use Galerkin projection or construct an interpolant to obtain an approximation to  $u$  in  $\mathcal{P}_{\Lambda_p}(\Gamma)$ . Let  $u_{\Lambda_p}$  denote the Galerkin projection of  $u$  onto the space  $\mathcal{P}_{\Lambda_p}(\Gamma)$ . Then we have the estimate

$$\|u - u_{\Lambda_p}\|_{L_\rho^2(\Gamma; H_0^1(D))} \leq C_a \min_{v \in H_0^1(D) \otimes \mathcal{P}_{\Lambda_p}(\Gamma)} \|u - v\|_{L_\rho^2(\Gamma; H_0^1(D))}$$

where  $C_a > 0$  depends on the coefficient  $a(x, \mathbf{y})$  and the bounds from assumption (A1). This estimate expresses optimality in the  $L_\rho^2(\Gamma)$  error of the Galerkin projection since  $C_a$  does not grow with  $\Lambda_p$ , and suggests that the Galerkin method is the best choice for approximating  $u$  in the space  $\mathcal{P}_{\Lambda_p}(\Gamma)$ . We can also define an interpolation operator  $\mathcal{I}_L^{\mathbf{m},g} : C^0(\Gamma) \rightarrow \mathcal{P}_{\Lambda_p}(\Gamma)$ , and then we have the estimate

$$\begin{aligned} \|u - \mathcal{I}_L^{\mathbf{m},g}[u]\|_{L_\rho^\infty(\Gamma; H_0^1(D))} &\leq (C_{\Lambda_L} + 1) \min_{v \in H_0^1(D) \otimes \mathcal{P}_{\Lambda_p}(\Gamma)} \|u - v\|_{L_\rho^\infty(\Gamma; H_0^1(D))} \\ &= (C_{\Lambda_L} + 1) \|u - u_{\Lambda_p}\|_{L_\rho^\infty(\Gamma; H_0^1(D))} \end{aligned} \quad (59)$$

where  $C_{\Lambda_L}$  is the Lebesgue constant of  $\mathcal{I}_L^{\mathbf{m},g}$ . A good interpolant will be one for which  $C_{\Lambda_L}$  grows moderately with  $\#\Lambda_L^{\mathbf{m},g}$ . For example, it is known (see [11, 16]) that for a one-dimensional Lagrange interpolation operator using a Clenshaw-Curtis rule, the Lebesgue constant is bounded by  $\frac{2}{\pi} \log(m-1) + 1$ , where  $m$  is the number of points. For the SC method, we define SDOF to be the total number of points needed to construct the approximation. From (59), if we only consider the number of SDOF needed to represent the solution, we expect the error for the Galerkin approximation to be much lower than the error in the interpolant. Indeed, this is reflected in our numerical results in Figures 7 and 9, and has been observed in previous comparisons [3, 13].

However, if we are willing to change the space  $\Lambda_p$ , e.g., adding more interpolation points to gain a more stable interpolant by changing  $\mathbf{m}$  or changing which points are included in the set  $\Lambda_L^{\mathbf{m},g}$  by changing  $g$ , it might be possible to obtain an approximation with lower complexity to reach a given tolerance, despite having to solve more systems. Therefore, to properly compare the work involved in constructing  $u_{\Lambda_p}$  and  $\mathcal{I}_L^{\mathbf{m},g}[u]$ , we consider the computational complexity of both methods, not in terms of SDOF, but in terms of floating point operations (FLOPs). For a chosen  $\Lambda_p$ , this reduces to studying the complexity of the system resulting from Galerkin projections and the stability properties of the interpolant  $\mathcal{I}_L^{\mathbf{m},g}$ .

Let  $\tilde{u}_{h,L}$  denote the numerical solution to the fully discrete approximation  $u_{h,L}$  obtained with the SCFEM from (53) found by the PCG, and observe that we have a similar splitting to (37) for the error in the approximation

$$\|u - \tilde{u}_{h,L}\|_{\mathcal{H}_\rho^2} \leq \underbrace{\|u - u_h\|_{\mathcal{H}_\rho^2}}_{\text{SC(I)}} + \underbrace{\|u_h - u_{h,L}\|_{\mathcal{H}_\rho^2}}_{\text{SC(II)}} + \underbrace{\|u_{h,L} - \tilde{u}_{h,L}\|_{\mathcal{H}_\rho^2}}_{\text{SC(III)}} \quad (60)$$

Note that unlike in the case of the SGFEM, the SCFEM does not require a further projection of the coefficient  $a(x, \mathbf{y})$ , so that we do not need to consider the error  $\|u - u^r\|_{\mathcal{H}_\rho^2}$  from (37). In addition, we do not need to

worry about well-posedness of the truncation as discussed in Remark 3.3. Similar to the complexity analysis for the SGFEM, we must choose  $h \leq h_{\max}$  and  $L \geq L_{\min}$  so that the errors  $\|u - u_h\|_{\mathcal{H}_\varepsilon^2}$  from the finite element discretization and  $\|u_h - u_{h,L}\|_{\mathcal{H}_\varepsilon^2}$  from the SC interpolation are both bounded by  $\varepsilon/3$ . From this, a minimum tolerance  $\tau_{\min}$  for the PCG solver can be derived and the maximum number of PCG iterations, with a zero initial guess, can be estimated [16]. In what follows, we present a result, whose proof can be found in [16, Theorem 4.7] that bounds the number of PCG iterations in the context of the work estimate (58). Using this estimate we can compare the cost in FLOPs for the SCFEM with the SGFEM results from Theorems 4.9 and 4.10 in the previous section.

**Theorem 5.1.** *Let  $u \in L_\varepsilon^2(\Gamma; H_0^1(D) \cap H^{s+1}(D))$  be the solution to (1). Then for  $\varepsilon > 0$  arbitrary, the work of finding  $\tilde{u}_{h,L}$ , the approximation to the fully discrete SC solution  $u_{h,L}$  from (53) found by PCG, denoted by  $W^{\text{PSC}}$ , can be bounded by*

$$W^{\text{PSC}} \leq 2C_7 \left( \frac{3C_{\text{FEM}}}{\varepsilon} \right)^{\frac{d}{s}} C_8 \left[ \log \left( \frac{3C_{\text{SC}}}{\varepsilon} \right) \right]^N \left[ C_9 + \frac{1}{\log 2} \log \log \left( \frac{3C_{\text{SC}}}{\varepsilon} \right) \right]^{N-1} \quad (61)$$

$$\times \frac{1}{\log \left( \frac{\sqrt{\bar{\kappa}}+1}{\sqrt{\bar{\kappa}}-1} \right)} \left\{ \log \left( \frac{C_{10}}{\varepsilon} \right) + C_{11} + 2N \log \log \left[ \frac{1}{rN} \log \left( \frac{3C_{\text{SC}}}{\varepsilon} \right) \right] \right\}.$$

Here  $C_{\text{FEM}}$  from Lemma 4.7,  $C_7$  from Theorem 4.9, and  $C_8, C_9, C_{10}, C_{11}, C_{\text{SC}}$ , and  $r$  from [16, Theorem 4.7] are positive constants independent of  $\varepsilon$ . Moreover, we define  $\bar{\kappa} = \sup_{\mathbf{y} \in \Gamma} \kappa(\mathbf{y})$  where  $\kappa(\mathbf{y})$  is the condition number of the preconditioned system  $\tilde{\mathbf{P}}_0^{-1} \mathbf{A}(\mathbf{y})$  with  $\mathbf{P}_0$  from (56).

Theorem 5.1 follows from the fact that  $N_{\text{iter}}^{\text{PSC}} = \sum_{k=1}^{M_L} N_{\text{iter}}^{\text{P}(k)} \leq N_{\text{zero}}$ , where  $N_{\text{zero}}$  is the number of iterations needed by the SCFEM with a PCG and a zero vector initial guess. Substituting the bound on  $N_{\text{zero}}$  shown in [16, Theorem 4.7] into the work estimate (58) and using  $J_{h_{\max}} = C_7(3C_{\text{FEM}}/\varepsilon)^{d/s}$  as in Theorem 4.9, puts the result in terms of FLOPs. Given Theorem 5.1 we see that the work of obtaining the fully discrete approximation with the SCFEM with the PCG method is asymptotically bounded by:

$$\underbrace{\mathcal{O} \left( \frac{1}{\varepsilon} \right)^{\frac{d}{s}}}_{(\text{SC.1})} \underbrace{\left[ \log \left( \frac{1}{\varepsilon} \right) \right]^N \left[ \log \log \left( \frac{1}{\varepsilon} \right) \right]^{N-1}}_{(\text{SC.2})} \underbrace{\left( \frac{\log \left( \frac{1}{\varepsilon} \right)}{\log \left( \frac{\sqrt{\bar{\kappa}}+1}{\sqrt{\bar{\kappa}}-1} \right)} \right)}_{(\text{SC.3})} \quad (62)$$

where (SC.1), (SC.2), and (SC.3) correspond to the work required by the finite element, SC interpolant, and PCG methods, respectively. Since the costs associated with the finite element discretization are the same for both methods, we focus only on the costs associated with the SG projection and the SC interpolation, coupled with the costs of the PCG method. In particular, as the work required by the SG approximation from (SG.2) of (47) has different bounds depending on whether the coefficient is a fixed order polynomial or is given by a total degree orthogonal expansion having order depending on  $\varepsilon$ , we now provide a comparison in both of these cases.

*Comparison in the affine and non-affine polynomial cases, e.g., Examples 2.1 and 2.2.* The terms (SG.2) from (47) and (SC.2) from (62) are asymptotic estimates of the number of coupled and decoupled finite element systems that must be solved by the SGFEM and SCFEM to construct the stochastic approximation, respectively. For the coefficients from Examples 2.1 and 2.2, (SG.2) from (47) for the SGFEM is  $\mathcal{O}([\log(1/\varepsilon)]^N)$ . Hence, our analysis shows that in these cases, the number of coupled finite element systems in the SGFEM matrix is smaller than the number of decoupled finite element systems required by the SCFEM by a factor of  $(\log \log(1/\varepsilon))^{N-1}$ . This difference is enough to suggest that, if the condition numbers of both the preconditioned coupled system from the SGFEM and the preconditioned individual systems from the SCFEM are of the same order, then the SGFEM will outperform the SCFEM in terms of minimum work required to obtain a fully discrete approximation. In fact, whenever  $a(x, \mathbf{y})$  is a general non-affine, polynomial coefficient, e.g., Example 2.2, having fixed order  $\bar{r} < \infty$ , the complexity of matrix-vector products

involving  $\mathbf{K}$  from (16) is approximately  $\mathcal{O}(J_h M_p M_{\bar{r}} \min\{2^{\bar{r}}, M_{\lceil \bar{r}/2 \rceil}\})$  from Remark 4.4. Hence, our analysis shows that when  $\mathcal{O}(M_{\bar{r}} \min\{2^{\bar{r}}, M_{\lceil \bar{r}/2 \rceil}\}) < \mathcal{O}((\log \log(1/\varepsilon))^{N-1})$ , which, in limit as  $\varepsilon \rightarrow 0$ , is always the case, and the condition numbers of both systems are of the same order, the SGFEM will outperform the SCFEM. However, in practical applications, it may require unrealistically small tolerance  $\varepsilon$  to see this when  $\bar{r}$  is large.

*Comparison in the non-affine, transcendental case, e.g., Example 2.3.* In the case that  $a(x, \mathbf{y})$  is a non-affine, transcendental function of the random parameters, e.g., Example 2.3, the estimate for (SG.2) is  $\mathcal{O}([\log(1/\varepsilon)]^{3N})$ . Here our analysis shows that the number of coupled finite element matrices present in the SG system  $\mathbf{K}_r$  from (22) dominates the number of decoupled finite element systems needed by the SCFEM by a factor of  $[\log(1/\varepsilon)]^{2N}$ . Note that, as in the case of the SGFEM, the term (SC.3) has a dependence on the condition numbers of the preconditioned finite element systems through the bound  $\bar{\kappa} = \sup_{\mathbf{y} \in \Gamma} \kappa(\mathbf{y})$ . For the unpreconditioned systems  $\mathbf{A}(\mathbf{y})$ , the condition numbers can be bounded by  $\kappa(\mathbf{A}(\mathbf{y})) \leq (C_\kappa/h)^2$  for every  $\mathbf{y} \in \Gamma$ , following from assumption (A1) and the quasi-uniformity of the mesh  $\mathcal{T}_h$ . However, if we use the exact inverse of  $\mathbf{P}_0$  when preconditioning the SCFEM systems, the condition numbers are bounded independent of  $h$  and  $\mathbf{y} \in \Gamma$ , since in this case  $\bar{\kappa}$  is independent of mesh size  $h$  and level  $L$ . Hence the work required by the PCG method when solving the SCFEM systems is dependent on  $\varepsilon$  only through the term  $\log(1/\varepsilon)$ . On the other hand, if we use the exact inverse of  $\mathbf{P}$  when preconditioning the SGFEM system, the condition number  $\tilde{\kappa}_r$  can be bounded by

$$\tilde{\kappa}_r \leq \frac{1 + \bar{\tau}_r}{1 - \underline{\tau}_r},$$

where  $\underline{\tau}_r$  and  $\bar{\tau}_r$  are defined in (49), hence depend on  $\varepsilon$  when  $r$  is chosen to satisfy  $r \geq r_{\min}$  from Lemma 4.8. Figure 5 plots the condition numbers of both the unpreconditioned matrix  $\mathbf{K}_r$  and the preconditioned matrix  $\mathbf{P}^{-1}\mathbf{K}_r$  with decreasing finite element mesh parameter  $h$  for the coefficient  $a(x, \mathbf{y})$  given in (67) from §6.3 with  $N = 4$ ,  $L_c = 1/2$ , and letting  $r = p$  with  $p$  increasing. There we see that the dependence on  $h$  has been removed by applying  $\mathbf{P}^{-1}$ , but as  $p$  increases, we see a corresponding increase in the condition number. Other preconditioners than  $\mathbf{P}$  may be used to reduce the dependence on  $r$ , e.g. [34], but then their associated costs must be accounted for in the work estimate (28) as well. However, even if the condition numbers of both the preconditioned coupled SG system and the preconditioned decoupled SC systems are of the same order, the additional work required to solve the coupled systems induced by the nonlinearity of the coefficient makes it difficult to see how the SGFEM can compete with the SCFEM.

## 6. Numerical examples

In this section, we provide illustrative numerical examples comparing the complexity of the SGFEM in the three cases of Examples 2.1, 2.2, and 2.3. We then compare these results with SCFEM and the results of the theoretical complexity comparison of the previous section. We solve the model problem (1), on the unit square  $D = [0, 1]^2$ . For a general coefficient  $a(x, \mathbf{y})$  we do not know the exact solution to (1). Hence we check the convergence against a “highly enriched” approximation, which we consider close enough to the exact one. To construct this “exact” solution  $u_{\text{ex}}(x, \mathbf{y})$ , we make use of the isotropic SCFEM based on Clenshaw-Curtis abscissas using the level  $L_{\text{ex}}$ . We approximate the computational error for the SGFEM with orders  $p = 0, 1, 2, \dots, p_{\max}$  and for the SCFEM with levels  $L = 0, 1, 2, \dots, L_{\max}$  as

$$\|\mathbb{E}[\varepsilon_{\text{SG}}]\|_{\ell^\infty} \approx \|\mathbb{E}[u_{\text{ex}} - \tilde{u}_{h,p}]\|_{\ell^\infty} \quad \text{and} \quad \|\mathbb{E}[\varepsilon_{\text{SC}}]\|_{\ell^\infty} \approx \|\mathbb{E}[u_{\text{ex}} - \tilde{u}_{h,L}]\|_{\ell^\infty}, \quad (63)$$

where  $\tilde{u}_{h,p}$  and  $\tilde{u}_{h,L}$  are the fully discrete approximations (13) and (53), respectively, found by the PCG method, described in §3 and §5. In §6.3, we measure  $\|\mathbb{E}[\varepsilon_{\text{SG}}]\|_{\ell^\infty} \approx \|\mathbb{E}[u_{\text{ex}} - \tilde{u}_{h,p}^r]\|_{\ell^\infty}$  where  $\tilde{u}_{h,p}^r$  denotes the solution of (22) with the projected coefficient  $a^r(x, \mathbf{y})$ .

As stated in §3.4 and §5.2, we use PCG with the mean-based preconditioner for SGFEM and the level-zero preconditioner for the SCFEM. Hence, we believe this puts both methods at a similar starting point for comparison, if not providing a slight advantage for the SGFEM. With these choices, the complexity results



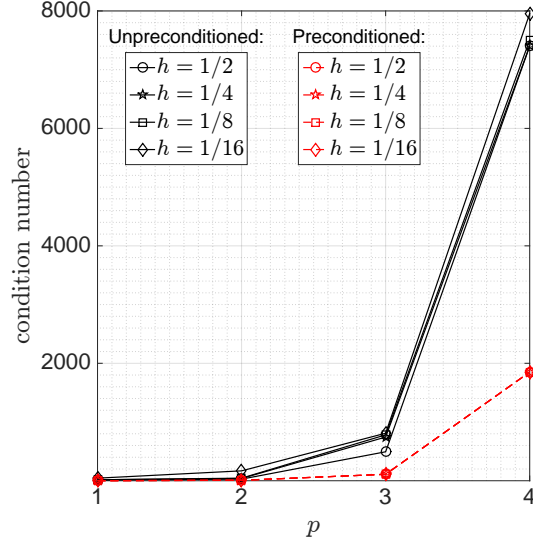


Figure 5: Condition numbers of both the unpreconditioned matrix  $\mathbf{K}_r$  and the preconditioned matrix  $\mathbf{P}^{-1}\mathbf{K}_r$  with decreasing finite element mesh parameter  $h$  and  $r = p$  for the coefficient (67) from §6.3 with  $N = 4$  and  $L_c = 1/2$ .

are presented in terms of the work estimates (28) and (58), respectively. The amount of work to reach a given error in PCG is also dependent on the tolerance used by the solver. If the tolerance is too small, we may see that the PCG method “over-resolves” the solution. To ensure that we do not over-resolve either solution, we set the tolerance of the solvers to be  $\|\mathbb{E}[\varepsilon_{SG}]\|_{\ell^\infty}/10$  and  $\|\mathbb{E}[\varepsilon_{SC}]\|_{\ell^\infty}/10$  respectively, where these quantities are first estimated for each order  $p$  and level  $L$  using a tolerance of  $1.0 \times 10^{-12}$ . In practice, we find that this does not affect the convergence results much.

In all three examples, we use the SG approximation constructed in terms of the orthonormal Legendre polynomials  $\{\Psi_{\mathbf{p}}\}_{\mathbf{p} \in \Lambda_p}$  for given index sets  $\Lambda_p$ . In the presentation of the results that follow, we use the following abbreviations. For the SGFEM, we use: “SG-TD” to denote the approximation in the total degree subspace  $\mathcal{P}_{\Lambda_p^{\text{TD}}}(\Gamma)$  with  $\Lambda_p^{\text{TD}}$  given in (12), and “SG-SM” to denote the approximation in the sparse Smolyak subspace  $\mathcal{P}_{\Lambda_p^{\text{SM}}}(\Gamma)$  with  $\Lambda_p^{\text{SM}}$  given in (12). For the SCFEM, we use: “SC-GL” and “SC-LJ” to denote the Smolyak approximation constructed on Gauss-Legendre abscissas and the Leja approximation constructed on Clenshaw-Curtis abscissas, both defined in terms of  $g_{\text{TD}}$  and  $\mathbf{m}$  given in (51) and (55), respectively, and “SC-CC” to denote the Smolyak approximation constructed on Clenshaw-Curtis abscissas with  $g_{\text{SM}}$  and  $\mathbf{m}$  given in (51) and (54).

### 6.1. Piecewise affine coefficients

One common example in engineering and the physical sciences is that of isotropic thermal diffusion problem with a stochastic conductivity coefficient. Consider a partitioning of  $D = [0, 1]^2$  into 8 circular inclusions arrayed about 1 square inclusion as in Figure 6. We present the following example from [3], where the coefficient was given by

$$a(x, \mathbf{y}) = b_0(x) + \sum_{n=1}^8 y_n \chi_n(x), \quad (64)$$

with  $b_0 = 1$ , and  $y_n \sim \mathcal{U}(-0.99, -0.2)$ . Here,  $\chi_n$  are indicator functions corresponding to the 8 circular inclusions of radius  $r = 0.13$ . In this example, we also set the forcing term to be

$$f(x) = 100\chi_F(x), \quad (65)$$

where  $F = [0.4, 0.6]^2$ , is the square inclusion centered in  $D$  with side length 0.2. Figure 6 shows the expected value of the solution to this problem. To solve (1) with the coefficient (64) and forcing function (65), we use a piecewise linear finite element basis in the deterministic space over a nonuniform mesh  $\mathcal{T}_h$ . Here, the nodes of  $\mathcal{T}_h$  are adapted to the geometry of our problem, that is, we fix the nodes that lie on the boundaries of the inclusions in our domain. From this fixed boundary data, we then use the `distmesh` MATLAB program [28] to generate a non-degenerate triangulation that adequately resolves the details of our subdomain geometry. We further specify the subsets of the total set of nodes that belong to each geometric inclusion, and to the boundaries of the inclusion, so that the interface conditions for the coefficient may be correctly applied. The final mesh consists of 10,604 elements, 5,377 total nodes, and 5,229 unknowns.

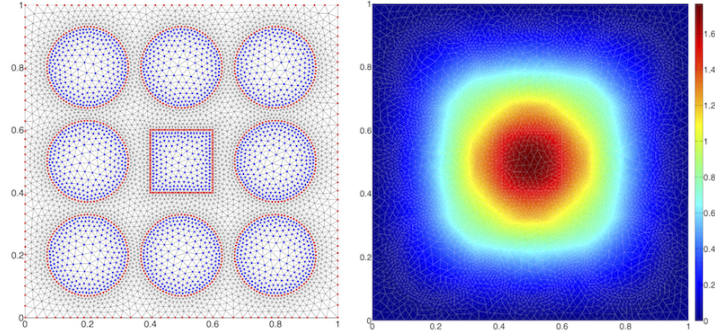


Figure 6: **Left:** a triangulation of the domain  $D$  with circular and square inclusions. Red nodes highlight the boundary of an inclusion or the domain  $D$ , blue nodes highlight nodes on the interior of an inclusion. **Right:** the expected value of the solution of (1) with stochastic conductivity coefficient (64).

The coefficient (64) is an example of a coefficient  $a(x, \mathbf{y})$  having affine dependence on the parameters, e.g., Example 2.1. Figure 7 displays the convergence of the stochastic Galerkin and collocation methods against the total number of SDOF. For the SGFEM we take the SDOF to be the cardinality of the set  $\Lambda_p$  used in constructing the fully discrete approximation  $u_{h,p}$  from (13) by solving (14), and for the SC method we take the SDOF to be the number of points  $\#\mathcal{G}_L^{m,g}$  corresponding to an index set  $\Lambda_L^{m,g}$  used in constructing the fully discrete approximation  $u_{h,L}$  from (53). From the discussion of §5.3, we expect to see that the approximation obtained with the SGFEM requires fewer SDOF than the SCFEM to achieve the same error, and this is indeed the observed result. For example, both the SG-TD and SC-LJ approximations require the same number of SDOF, but the error of the SC-LJ approximation is much higher. This, of course, is a consequence of the estimate (59), where the errors of the SC approximations are bounded above by their respective Lebesgue constants against the best-approximation error in the space  $\mathcal{P}_{\Lambda_L^{m,g}}(\Gamma)$ .

Figure 7 also displays the convergence of both methods in terms of error versus the total computational cost of solving the system with the work estimates of (28) and (58), respectively. Here, we compute the error in  $\|\mathbb{E}[\varepsilon_{SG}]\|_{\ell^\infty}$  and  $\|\mathbb{E}[\varepsilon_{SC}]\|_{\ell^\infty}$  as given in (60) and measure the cost as the number of  $\mathcal{O}(J_h)$  matrix vector products required by both methods which are explicitly counted as

$$N_{\text{iter}}^{\text{pSG}} * \left( M_p + \sum_{\mathbf{r} \in \Lambda_r} \text{nnz}(\mathbf{G}_{\mathbf{r}}) \right) = N_{\text{iter}}^{\text{pSG}} * (M_p + \mathcal{M}(p, r))$$

in the code for the SGFEM and  $2 * \sum_{k=1}^{M_L} N_{\text{iter}}^{\text{p}(k)}$  in the code for the SCFEM. Our analysis shows the work corresponding to the SG discretization for SG-TD is asymptotically bounded by  $\mathcal{O}([\log(1/\varepsilon)]^N)$  while the analysis from [16] shows that the work corresponding to the SC discretization for SC-CC is asymptotically bounded by  $\mathcal{O}([\log(1/\varepsilon)]^N [\log \log(1/\varepsilon)]^{N-1})$ . This closely matches the results of the numerical experiments in Figure 7, where it can be seen that for polynomial order  $p \geq 2$ , the SG-TD approximation yields the best results with the least computational cost for this problem.

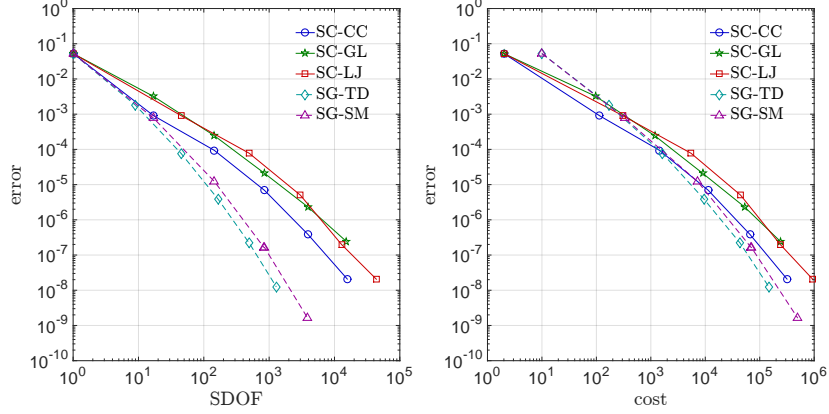


Figure 7: **Left:** Error versus SDOF in solving problem (1) with coefficient (64) and forcing (65). **Right:** Error versus computational cost with the work estimates given in (28) and (58) based on total number of matrix-vector products used by the CG method.

### 6.2. Polynomial coefficients

The next example we present is that of a polynomial function of the random parameters  $\mathbf{y}$ , e.g. the coefficient from Example 2.2. We consider the following function

$$a(x, \mathbf{y}) = 5 + \sum_{|\mathbf{r}| \leq \bar{r}} e^{-1.5|\mathbf{r}|} \varsigma_{\mathbf{r}}(x) \mathbf{y}^{\mathbf{r}}, \quad \varsigma_{\mathbf{r}}(x) = \begin{cases} \sin(|\mathbf{r}|\pi x_1) \cos(|\mathbf{r}|\pi x_2) & \text{if } |\mathbf{r}| \text{ is even,} \\ \cos(|\mathbf{r}|\pi x_1) \sin(|\mathbf{r}|\pi x_2) & \text{if } |\mathbf{r}| \text{ is odd,} \end{cases} \quad (66)$$

with  $y_n \sim \mathcal{U}(-1, 1)$  for all  $n = 1, \dots, N$  and forcing term  $f(x) = 1 \forall x \in \bar{D}$ . For the results that follow we fix  $N = 4$  and study the convergence of the SGFEM and SCFEM in the cases  $\bar{r} = 1, 3, 7$  in (66). As in §6.1, we set the finite element space for the spatial discretization to be the span of piecewise linear polynomials, but here we use a uniform triangulation of  $D$  with 4,934 elements and 2,340 spatial unknowns.

Figure 8 displays the convergence of the SGFEM and SCFEM in terms of error versus the total computational cost of solving the system with the work estimates of (28) and (58). Here, we compute the error in  $\|\mathbb{E}[\varepsilon_{\text{SG}}]\|_{\ell^\infty}$  and  $\|\mathbb{E}[\varepsilon_{\text{SC}}]\|_{\ell^\infty}$  as given in (63). As we increase the order  $\bar{r}$  in (66), we see that the work for the SGFEM increases, corresponding to the decreasing sparsity of the matrix  $\mathbf{K}$  from (16). Here, the work of matrix-vector multiplications with  $\mathbf{K}$  are of the order  $\mathcal{O}(J_h M_p M_{\bar{r}} \min\{2^{\bar{r}}, M_{\lceil \bar{r}/2 \rceil}\})$ , where  $M_{\bar{r}} \min\{2^{\bar{r}}, M_{\lceil \bar{r}/2 \rceil}\}$  is a large constant that grows rapidly with  $\bar{r}$ . As a result, we see that for  $\bar{r} = 1$ , the SGFEM outperforms the other methods for  $p \geq 4$ . However, for  $\bar{r} = 3, 7$ , the extra work of the matrix-vector multiplications of the coupled SG system dominates the overall convergence. We also observe that the convergence rate of the SGFEM does not change in these cases, as discussed in the comparison in §5.3.

### 6.3. Transcendental coefficients

The next example we present is that of a random coefficient defined in terms of the truncated Karhunen-Loève expansion of the function  $\log(a(x, \mathbf{y}) - a_{\min})$ , for  $a_{\min} > 0$ . This example represents a commonly used transcendental function of the physical and random parameters, e.g., Example 2.3, and is often presented in the context of enforcing the positivity of  $a(x, \mathbf{y})$  required in assumption (A1). Coefficients of this type are commonly found in groundwater flow models. For these models, the permeability can exhibit large variance within each layer of sediment, and as a result are better represented on a logarithmic scale. We recall the problem of solving (1) with a coefficient having one-dimensional (layered) spatial dependence and

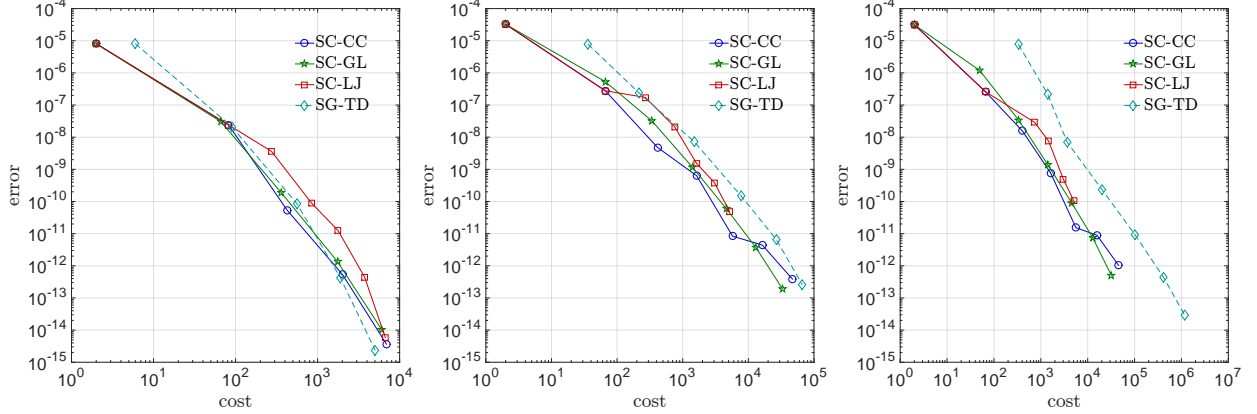


Figure 8: Error versus cost for solving problem (1) with coefficient (66) having  $\bar{r} = 1$  (left),  $\bar{r} = 3$  (middle), and  $\bar{r} = 7$  (right), with forcing  $f(x) = 1$ . The cost, given in (28) and (58), is based on total number of matrix-vector products used by the PCG method.

a deterministic load  $f(x_1, x_2, \omega) = 2 \cos(x_1) \sin(x_2)$  as studied in [25, 26], where  $a(x, \mathbf{y})$  was given by

$$\log(a(x, \omega) - 0.5) = 1 + y_1(\omega) \left( \frac{\sqrt{\pi}L}{2} \right)^{1/2} + \sum_{n=2}^N \zeta_n \varphi_n(x) y_n(\omega), \quad (67)$$

$$\zeta_n := (\sqrt{\pi}L)^{1/2} \exp \left( -(\lfloor \frac{n}{2} \rfloor \pi L)^2 / 8 \right), \quad \text{for } n > 1, \quad \varphi_n(x) := \begin{cases} \sin \left( \lfloor \frac{n}{2} \rfloor \pi x_1 / L_p \right), & \text{if } n \text{ is even,} \\ \cos \left( \lfloor \frac{n}{2} \rfloor \pi x_1 / L_p \right), & \text{if } n \text{ is odd.} \end{cases}$$

Here,  $\{y_n(\omega)\}_{n=1}^{\infty}$  are independent random variables uniformly distributed in  $[-\sqrt{3}, \sqrt{3}]$  with zero mean and unit variance. For  $x_1 \in [0, b]$ , let  $L_c$  be a desired physical correlation length for the random field  $a(x, \mathbf{y})$ , chosen so that the random variables  $a(x_1, \omega)$  and  $a(x'_1, \omega)$  become essentially uncorrelated for  $|x_1 - x'_1| \gg L_c$ . Also, let  $L_p = \max\{b, 2L_c\}$  and  $L = L_c/L_p$ . Expression (67) represents a possible truncation of a one-dimensional random field with stationary covariance,

$$\text{cov}[\log(a - 0.5)](x_1, x_2) = \exp \left( \frac{-(x_1 - x_2)^2}{L_c^2} \right).$$

Direct integration with the coefficient  $a(x, \mathbf{y})$  from (67) yields a fully-block dense linear system  $\mathbf{K}$  from (16) that is computationally infeasible to solve [14, 23, 34, 35]. The purpose of this example is to highlight the difficulties of obtaining a fully discrete approximation with the SGFEM in this case.

As in the previous example in §6.2, we set the finite element space for the spatial discretization to be the span of piecewise linear polynomials and use a uniform triangulation of  $D$  with 4,934 elements and 2,340 spatial unknowns. For the results that follow, we fix the truncation length  $N = 9$  and correlation length  $L_c = 1/64$  in (67). To maintain sparsity of the SG system, we use the strategy of projecting the coefficient  $a(x, \mathbf{y})$  from (67) onto the space  $\mathcal{P}_{\Lambda_r}(\Gamma)$ , as in (18), where  $\Lambda_r = \Lambda_r^{\text{TD}}$  for the SG-TD approximation, and  $\Lambda_r = \Lambda_r^{\text{SM}}$  for the SG-SM approximation, obtaining the matrix  $\mathbf{K}_r$  from (22). We then increase  $r$  while  $p$  is fixed until the error in the solution stagnates, in practice finding that, for this problem,  $r = p$  is sufficient to guarantee the error of the projection does not exceed that of the solution, while maintaining sparsity of the linear system.

Figure 9 compares the error versus SDOFs. There we see that for order  $p \geq 3$ , the SG-TD approximation provides the best approximation with respect to SDOFs. As discussed in §5.3, this is to be expected since the computational complexity of solving the coupled and decoupled systems is not taken into account. Figure 9 also displays the convergence in error versus the total computational cost of solving the system with the work estimates of (28) and (58). Here however, the results show that the SGFEM requires significantly

| SC-CC Level | SC-CC Error              | Mat-vec cost of SC-CC | SG-TD Order | SG-TD Error              | Mat-vec cost of SG-TD |
|-------------|--------------------------|-----------------------|-------------|--------------------------|-----------------------|
| 0           | $1.3626 \times 10^{-4}$  | 2                     | 0           | $1.3626 \times 10^{-4}$  | 4                     |
| 1           | $2.8884 \times 10^{-6}$  | 218                   | 1           | $3.9444 \times 10^{-5}$  | 152                   |
| 2           | $6.3652 \times 10^{-8}$  | 3,398                 | 2           | $6.1427 \times 10^{-7}$  | 10,710                |
| 3           | $3.6021 \times 10^{-9}$  | 28,638                | 3           | $2.8851 \times 10^{-8}$  | 213,010               |
| 4           | $1.4794 \times 10^{-10}$ | 178,894               | 4           | $4.9210 \times 10^{-10}$ | 4,579,575             |
| 5           | $2.2869 \times 10^{-12}$ | 944,220               | 5           | $8.9123 \times 10^{-12}$ | 49,089,051            |

Table 1: Comparison of cost in matrix-vector products for solving problem (1) with coefficient (67) and forcing  $f(x_1, x_2, \omega) = \cos(x_1) \sin(x_2)$  using the SC-CC and SG-TD approximations, with the strategy of picking the CG tolerance to be  $\|\mathbb{E}[u_{ex} - \tilde{u}_{h,p}^*]\|_{\ell^\infty}/10$  for the SGFEM and  $\|\mathbb{E}[u_{ex} - \tilde{u}_{h,L}]\|_{\ell^\infty}/10$  for the SC method. Cost in matrix-vector products for SG-TD method is given by (28) and for SC-CC is given by (58) normalized by the cost of a finite element matrix vector product.

more work to obtain the same error than the SCFEM. We also observe the change in rate discussed in §5.3 in this case, as the work required to solve (1) with the coefficient  $a(x, \mathbf{y})$  from (67) now depends on the order  $r$  of the projection used in the approximation of  $a(x, \mathbf{y})$ .

For the TD-SG approximation, this can be seen as a consequence of the fact that when  $r = p$ , the cost of solving (22) with the PCG method is of the order  $\mathcal{O}(J_h M_{[p/2]}^3 N_{\text{iter}}^{\text{SG}})$ , growing much more rapidly than the cost in the affine and polynomial coefficient cases, e.g., Examples 2.1 and 2.2, as we increase the order  $p$ . Table 1 shows the amount of work required to achieve an error on the order of  $10^{-k}$  for some values of  $k$  in terms of the total number of matrix-vector products required by both the SC-CC and SG-TD approximations.

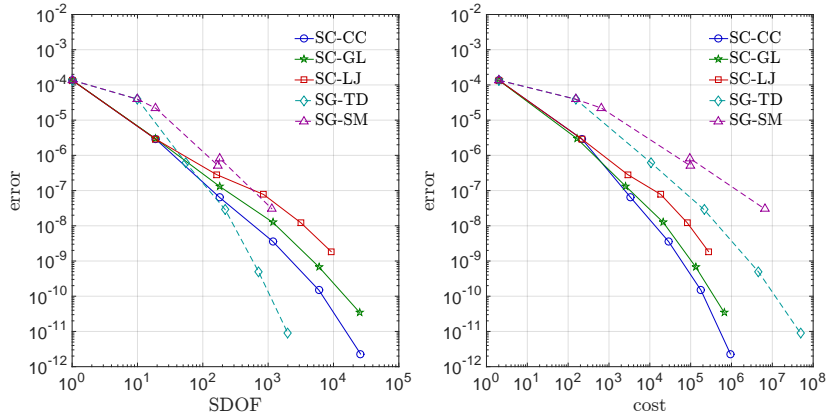


Figure 9: **Left:** Error versus SDOFs in solving problem (1) with coefficient (67) and forcing  $f(x_1, x_2, \omega) = \cos(x_1) \sin(x_2)$ . **Right:** Error versus cost with the work estimates given in (28) and (58).

## 7. Conclusions

In this work, we presented explicit cost bounds for applying the SGFEM to the solution of an elliptic PDE having both affine and non-affine random coefficients. To this end, we have conducted a rigorous counting argument for the sparsity of the linear system that results from the SG discretization with a global orthogonal basis defined on an isotropic total degree index set. Our analysis shows that when the coefficient is an affine or non-affine function of the random variables having fixed polynomial order, the computational cost of solving the coupled SG system grows linearly with the dimension of the polynomial subspace. In

these cases, the results only differ by a constant depending on the polynomial order of the random coefficient and the dimension of the parameter domain.

On the other hand, when the coefficient is a non-affine, transcendental function of the random variables requiring an additional orthogonal expansion, our analysis shows that the computational complexity, no longer grows linearly with the polynomial subspace dimension. For such coefficients, we are able to provide bounds on the complexity that depend on the truncation order of the coefficient. These estimates imply that a truncation of the expansion should be used, when possible, though attention must be paid to the well-posedness of the resulting PDE.

The analysis conducted herein motivates the study of the total computational complexity of obtaining fully discrete approximations with such methods. We have seen that, despite the fact that the SG method yields an approximation that is optimal in the  $L^2$  sense for a given polynomial subspace, the associated computational costs of obtaining SG approximations are not optimal for all problems. Moreover, we have observed, both through theoretical comparisons and numerical examples, that changing the underlying polynomial subspace and method used for obtaining the fully discrete approximation can often yield a solution that requires far less work to obtain, but has the same error.

### Acknowledgements

The first author would like to acknowledge Dr. Miroslav Stoyanov for his insightful comments and assistance in producing the stochastic collocation results with the TASMANIAN package [1].

### References

- [1] *Toolkit for Adaptive Stochastic Modeling and Non-Intrusive Approximation*. <http://tasmanian.ornl.gov/>.
- [2] K. ATKINSON AND W. HAN, *Theoretical Numerical Analysis: A Functional Analysis Framework*, vol. 39 of Texts in Applied Mathematics, Springer, New York, NY, 3rd ed., 2005.
- [3] J. BÄCK, F. NOBILE, L. TAMELLINI, AND R. TEMPONE, *Stochastic Spectral Galerkin and Collocation Methods for PDEs with Random Coefficients: A Numerical Comparison*, in Spectral and High Order Methods for Partial Differential Equations, J. S. Hesthaven and E. M. Rønquist, eds., vol. 76 of Lecture Notes in Computational Science and Engineering, Springer Berlin Heidelberg, 2011, pp. 43–62.
- [4] A. CHKIFA, A. COHEN, AND C. SCHWAB, *Breaking the curse of dimensionality in sparse polynomial approximation of parametric PDEs*, Journal de Mathématiques Pures et Appliquées, 103 (2014), pp. 400–428.
- [5] C. W. CLENSHAW AND A. R. CURTIS, *A method for numerical integration on an automatic computer*, Numerische Mathematik, 2 (1960), pp. 197–205.
- [6] A. COHEN, A. CHKIFA, R. DEVORE, AND C. SCHWAB, *Sparse adaptive Taler approximation algorithms for parametric and stochastic elliptic PDEs*, ESAIM: Mathematical Modelling and Numerical Analysis, (2012), pp. 1–27.
- [7] A. COHEN, R. DEVORE, AND C. SCHWAB, *Convergence Rates of Best N-term Galerkin Approximations for a Class of Elliptic sPDEs*, Foundations of Computational Mathematics, 10 (2010), pp. 615–646.
- [8] A. COHEN, R. DEVORE, AND C. SCHWAB, *Analytic regularity and polynomial approximation of parametric and stochastic elliptic PDEs*, Analysis and Applications, 09 (2011), pp. 11–47.
- [9] P. J. DAVIS, *Interpolation and Approximation*, Dover, 1975.
- [10] S. DE MARCHI, *On Leja sequences: Some results and applications*, Applied Mathematics and Computation, 152 (2004), pp. 621–647.
- [11] V. K. DZJADYK AND V. V. IVANOV, *On asymptotics and estimates for the uniform norms of the Lagrange interpolation polynomials corresponding to the Chebyshev nodal points*, Analysis Mathematica, 9 (1983), pp. 85–97.
- [12] M. EIERMANN, O. G. ERNST, AND E. ULLMANN, *Computational aspects of the stochastic finite element method*, in Computing and Visualization in Science, vol. 10, 2007, pp. 3–15.
- [13] H. C. ELMAN, C. W. MILLER, E. T. PHIPPS, AND R. S. TUMINARO, *Assessment of collocation and Galerkin approaches to linear diffusion equations with random data*, International Journal for Uncertainty Quantification, 1 (2011), pp. 19–33.
- [14] O. G. ERNST AND E. ULLMANN, *Stochastic Galerkin Matrices*, SIAM Journal on Matrix Analysis and Applications, 31 (2010), pp. 1848–1872.
- [15] G. FISHMAN, *Monte Carlo: Concepts, Algorithms, and Applications*, Springer Series in Operations Research and Financial Engineering, Springer, 1996.
- [16] D. GALINDO, P. JANTSCH, C. G. WEBSTER, AND G. ZHANG, *Accelerating stochastic collocation methods for PDEs with random input data*, Tech. Rep. ORNL/TM-2015/219, Oak Ridge National Laboratory, 2015.
- [17] M. GUNZBURGER AND C. G. WEBSTER, *Uncertainty quantification for partial differential equations with stochastic coefficients*, The Mathematical Intelligencer, (2014). To appear.
- [18] M. GUNZBURGER, C. G. WEBSTER, AND G. ZHANG, *An adaptive wavelet stochastic collocation method for irregular solutions of partial differential equations with random input data*, in Sparse Grids and Applications - Munich 2012, vol. 97 of Lecture Notes in Computational Science and Engineering, Springer International Publishing, 2014, pp. 137–170.

- [19] M. D. GUNZBURGER, C. G. WEBSTER, AND G. ZHANG, *Stochastic finite element methods for partial differential equations with random input data*, Acta Numerica, 23 (2014), pp. 521–650.
- [20] C. JOHNSON, *Numerical solution of partial differential equations by the finite element method*, Courier Dover Publications, 2012.
- [21] O. P. LE MAÎTRE AND O. M. KNIO, *Spectral Methods for Uncertainty Quantification: With Applications to Computational Fluid Dynamics*, Scientific Computation, Springer, 2010.
- [22] M. LOÈVE, *Probability Theory*, no. v. 2 in Graduate Texts in Mathematics, Springer, 1978.
- [23] H. G. MATTHIES AND A. KEESE, *Galerkin methods for linear and nonlinear elliptic stochastic partial differential equations*, Computer Methods in Applied Mechanics and Engineering, 194 (2005), pp. 1295–1331.
- [24] F. NOBILE AND R. TEMPONE, *Analysis and implementation issues for the numerical approximation of parabolic equations with random coefficients*, Online, 80 (2009), pp. 979–1006.
- [25] F. NOBILE, R. TEMPONE, AND C. G. WEBSTER, *A Sparse Grid Stochastic Collocation Method for Partial Differential Equations with Random Input Data*, SIAM Journal on Numerical Analysis, 46 (2008), pp. 2309–2345.
- [26] ———, *An Anisotropic Sparse Grid Stochastic Collocation Method for Partial Differential Equations with Random Input Data*, SIAM Journal on Numerical Analysis, 46 (2008), pp. 2411–2442.
- [27] M. F. PELLISSETTI AND R. G. GHANEM, *Iterative Solution of Systems of Linear Equations Arising in the Context of Stochastic Finite Elements*, Advances in Engineering Software, 31 (2000), pp. 607–616.
- [28] P.-O. PERSSON AND G. STRANG, *A Simple Mesh Generator in MATLAB*, SIAM Review, 46 (2004), pp. 329–345.
- [29] C. E. POWELL AND H. C. ELMAN, *Block-diagonal preconditioning for spectral stochastic finite-element systems*, IMA Journal of Numerical Analysis, 29 (2009), pp. 350–375.
- [30] S. SMOLYAK, *Quadrature and interpolation formulas for tensor products of certain classes of functions*, Dokl. Akad. Nauk SSSR, (1963), pp. 4:240–243.
- [31] G. SZEGÖ, *Orthogonal polynomials*, vol. XXIII, Amer. Math. Soc., 4 ed., 1975.
- [32] R. A. TODOR AND C. SCHWAB, *Convergence rates for sparse chaos approximations of elliptic problems with stochastic coefficients*, IMA Journal of Numerical Analysis, 27 (2007), pp. 232–261.
- [33] H. TRAN, C. G. WEBSTER, AND G. ZHANG, *Analysis of quasi-optimal polynomial approximations for parameterized PDEs with deterministic and stochastic coefficients*, Tech. Rep. ORNL/TM-2014/468, Oak Ridge National Laboratory, 2014. Submitted.
- [34] E. ULLMANN, *A Kronecker Product Preconditioner for Stochastic Galerkin Finite Element Discretizations*, SIAM Journal on Scientific Computing, 32 (2010), pp. 923–946.
- [35] E. ULLMANN, H. C. ELMAN, AND O. G. ERNST, *Efficient Iterative Solvers for Stochastic Galerkin Discretizations of Log-Transformed Random Diffusion Problems*, SIAM Journal on Scientific Computing, 34 (2012), pp. A659–A682.
- [36] N. WIENER, *The Homogeneous Chaos*, American Journal of Mathematics, 60 (1938), pp. pp. 897–936.
- [37] D. XIU AND G. E. KARNIADAKIS, *The Wiener–Askey Polynomial Chaos for Stochastic Differential Equations*, SIAM J. Sci. Comput., 24 (2002), pp. 619–644.

## 8. Appendix

PROOF (OF COROLLARY 4.3). When  $N = 1$  we denote  $\mathbf{r} = r \in \mathbb{N}_0$  and note that from Theorem 4.1,

$$c(r, \ell) = \begin{cases} \#\mathbf{S}(r, \ell) & r \text{ even, } \ell = r/2 \\ 2\#\mathbf{S}(r, \ell) & \text{otherwise,} \end{cases}$$

with  $\mathbf{S}(r, \ell) = \{s \in \mathbb{N}_0 : s = \ell, s \leq r\} = \{\ell\}$  for  $\ell \leq r$  and  $\emptyset$  otherwise. We distinguish in cases:

- Case  $r = 2k$ ,  $k \in \mathbb{N}_0$ ,

1. when  $0 \leq r \leq p$ ,

$$\begin{aligned} \text{nnz}(\mathbf{G}_r) &= \sum_{\ell=\lceil r/2 \rceil}^r c(r, \ell) \binom{1+p-\ell}{p-\ell} = \binom{1+p-k}{p-k} + \sum_{\ell=k+1}^{2k} 2 \binom{1+p-\ell}{p-\ell} \\ &= (1+p-k) - k(3k-2p-1) \\ &= 1+p-4k^2+2kp+k^2 \\ &= (p-2k+1)(2k+1)+k^2 \\ &= (p-r+1)(r+1)+k^2. \end{aligned}$$

2. when  $p+1 \leq r \leq 2p$ , we have  $\frac{p+1}{2} \leq k \leq p$ , so

$$\begin{aligned} \text{nnz}(\mathbf{G}_r) &= \sum_{\ell=\lceil r/2 \rceil}^r c(r, \ell) \binom{1+p-\ell}{p-\ell} = \binom{1+p-k}{p-k} + \sum_{\ell=k+1}^{2k} 2 \binom{1+p-\ell}{p-\ell} \\ &= (1+p-k) + \sum_{\ell=k+1}^p 2 \binom{1+p-\ell}{p-\ell} \\ &= (1+p-k) + (p-k)(p-k+1) \\ &= (1+p-k)^2. \end{aligned}$$

3. when  $r > 2p$ , then  $k > p$ , so

$$\text{nnz}(\mathbf{G}_r) = \sum_{\ell=\lceil r/2 \rceil}^r c(r, \ell) \binom{1+p-\ell}{p-\ell} = \binom{1+p-k}{p-k} + \sum_{\ell=k+1}^{2k} 2 \binom{1+p-\ell}{p-\ell} = 0,$$

since  $p-k < 0$  and  $l > k \Rightarrow p-\ell < p-k < 0$ .

• Case  $r = 2k+1$ ,  $k \in \mathbb{N}_0$ ,

1. when  $0 \leq r \leq p$ , then  $\lceil r/2 \rceil = \lceil (2k+1)/2 \rceil = \lceil k+1/2 \rceil = k+1$ , so

$$\begin{aligned} \text{nnz}(\mathbf{G}_r) &= \sum_{\ell=\lceil r/2 \rceil}^r c(r, \ell) \binom{1+p-\ell}{p-\ell} = 2 \sum_{\ell=k+1}^{2k+1} \binom{1+p-\ell}{p-\ell} \\ &= -(1+k)(3k-2p) \\ &= -4k+2p-4k^2+2kp+k^2+k \\ &= -2k(2k+2)+p(2k+2)+k^2+k \\ &= (p-2k)(2k+2)+k^2+k \\ &= (p-r+1)(r+1)+k^2+k. \end{aligned}$$

2. when  $p+1 \leq r \leq 2p$ , then  $p/2 \leq k \leq p-1/2$ , so

$$\begin{aligned} \text{nnz}(\mathbf{G}_r) &= \sum_{\ell=\lceil r/2 \rceil}^r c(r, \ell) \binom{1+p-\ell}{p-\ell} = 2 \sum_{\ell=k+1}^{2k+1} \binom{1+p-\ell}{p-\ell} = 2 \sum_{\ell=k+1}^p \binom{1+p-\ell}{p-\ell} \\ &= (p-k)(p-k+1). \end{aligned}$$

3. when  $r > 2p$ , then  $k > p-1/2$ , so

$$\text{nnz}(\mathbf{G}_r) = \sum_{\ell=\lceil r/2 \rceil}^r c(r, \ell) \binom{1+p-\ell}{p-\ell} = \sum_{\ell=k+1}^{2k+1} c(r, \ell) \binom{1+p-\ell}{p-\ell} = 0,$$

since  $k > p-1/2 \Rightarrow p-\ell \leq p-(k+1) = p-k-1 < p-(p-1/2)-1 = -1/2$ . □